

1 **Tonic dopamine and biases in value learning linked through a**
2 **biologically inspired reinforcement learning model**

3

4

5 Sandra Romero Pinto^{1,2}, Naoshige Uchida¹

6

7 Affiliations:

8 ¹Department of Molecular and Cellular Biology, Center for Brain Science, Harvard University,

9 Cambridge, MA 02138, USA

10 ²Program in Speech and Hearing Bioscience and Technology, Division of Medical Sciences,

11 Harvard Medical School, Boston, MA 02115, USA

12

13 Correspondence: uchida@mcb.harvard.edu (N.U); sromeropinto@g.harvard.edu (S.R.P)

14

15

16

17 **Abstract**

18

19 A hallmark of various psychiatric disorders is biased future predictions. Here we examined the
20 mechanisms for biased value learning using reinforcement learning models incorporating recent
21 findings on synaptic plasticity and opponent circuit mechanisms in the basal ganglia. We show
22 that variations in tonic dopamine can alter the balance between learning from positive and
23 negative reward prediction errors, leading to biased value predictions. This bias arises from the
24 sigmoidal shapes of the dose-occupancy curves and distinct affinities of D1- and D2-type
25 dopamine receptors: changes in tonic dopamine differentially alters the slope of the dose-
26 occupancy curves of these receptors, thus sensitivities, at baseline dopamine concentrations. We
27 show that this mechanism can explain biased value learning in both mice and humans and may
28 also contribute to symptoms observed in psychiatric disorders. Our model provides a foundation
29 for understanding the basal ganglia circuit and underscores the significance of tonic dopamine in
30 modulating learning processes.

31

32

33 **Introduction**

34

35 Our ability to predict the outcomes of our actions is crucial in selecting and motivating
36 appropriate actions. Systematic biases in future predictions or expectations, however, can lead to
37 maladaptive behaviors, such as those observed in patients with various psychiatric disorders¹⁻⁴.
38 For example, overly negative or pessimistic predictions can contribute to major depression^{1,5},
39 whereas excessively positive or optimistic predictions may be associated with pathological
40 gambling, addiction, and mania^{3,4,6-8}. Despite the importance of understanding the causes of
41 biased future predictions, the biological mechanisms underlying them remain poorly understood.

42

43 Our future expectations and decisions are shaped by experiences of positive and negative events.
44 The process of learning from outcomes has been modeled using reinforcement learning (RL)
45 models⁹⁻¹², where value predictions are updated based on reward prediction errors (RPEs), that is
46 the discrepancy between received and expected outcomes. In addition to its role in learning,

47 recent studies have indicated the importance of RPEs in mood; these studies have suggested that
48 mood depends not on the absolute goodness of outcomes, but rather on the recent history of
49 RPEs^{13,14}.

50

51 In the brain, dopamine is thought to be a key regulator in this process of learning from positive
52 and negative outcomes. The dynamics of dopamine are often categorized into two modes: tonic
53 and phasic. Tonic dopamine refers to “baseline” dopamine that operates on a long timescale such
54 as tens of seconds or minutes, while phasic activity refers to transient changes that occur at a
55 much shorter, sub-second timescale, often triggered by external stimuli^{15–18}. A significant body
56 of evidence has shown that phasic responses of dopamine neurons convey RPEs and drive
57 learning of values and actions^{17–20}. On the other hand, changes in tonic dopamine might also
58 modulate value learning, yet whether and how the level of tonic dopamine modulates learning
59 remain poorly understood.

60

61 Previous studies have reported that patients with psychiatric disorders exhibit biased learning
62 from positive versus negative outcomes. For one, some studies have shown that patients with
63 major depression have a reduced sensitivity in learning from rewarding events, while their ability
64 to learn from negative events remains relatively intact^{1,5,21}. Similarly, patients with Parkinson’s
65 disease are better at learning from negative than positive outcomes^{22,23}. Analysis of these patients
66 using RL models has suggested that biases in learning can be explained by alterations in specific
67 parameters in RL models, such as the learning rate parameters or the sensitivity to positive and
68 negative outcomes. For example, some studies have suggested that anhedonia in major
69 depressive disorder may correspond to a reduced learning rate from positive compared to
70 negative outcomes¹.

71

72 Mechanistically, some of these changes in RL parameters can be linked to altered functions of
73 dopamine. First, it has been shown that dopamine synthesis capacity, an approximate indicator of
74 baseline dopamine levels, in the striatum, as measured using positron emission tomography

75 (PET), correlates with learning rate parameters²⁴. Second, dopamine medications can change the
76 balance between learning from positive and negative outcomes^{22,24,25}. Third, responses to
77 positive outcomes in the nucleus accumbens (NAc), as measured based on blood oxygenation-
78 dependent (BOLD) signals, are reduced in patients with psychiatric disorders such as
79 depression²⁶⁻²⁹. These observations point to important roles of reinforcement learning processes
80 and dopamine in regulating value learning. However, the parameters in RL models remain an
81 abstract entity, and biological processes underlying changes in these parameters are still largely
82 unknown.

83
84 One limitation in most RL models used in previous studies is that they do not reflect key neural
85 circuit architectures in the brain (but see ³⁰⁻³²) nor recent findings on intracellular signaling and
86 plasticity rules that can constrain how dopamine functions in biological circuits³³⁻³⁵.
87 Incorporating these key biological factors may lead to better understanding of how changes in
88 RL parameters may arise in psychiatric disorders. Furthermore, recent studies have found that
89 the activity of dopamine neurons is consistent with a novel RL algorithm called distributional
90 RL³⁶⁻³⁸. Distributional RL takes into account the diversity in dopamine signals, and a population
91 of dopamine neurons together encodes the entire distribution of rewards, not just the average.
92 Although distributional RL has shown to be efficient in solving various RL problems in artificial
93 intelligence^{37,39}, how distributional RL can be implemented in biological neural circuits and how
94 distributional RL relates to biased value learning remain to be examined.

95
96 In this study, we sought to identify potential mechanisms that cause biased value predictions
97 using biologically inspired RL models. To this goal, we first construct an RL model that
98 incorporates recent biological findings, such as intracellular signaling and synaptic plasticity
99 rules as well as the basic circuit architecture in the brain³². Based on this model, we propose two
100 potential biological mechanisms that can cause optimistic or pessimistic biases in value
101 predictions. We will then show that some existing data can be explained by one of these models.

102 Finally, we will show how our model can provide an account of how biases in value predictions
103 arise in psychiatric disorders.

104

105

106 **Results**

107

108 **Basic reinforcement learning algorithms**

109 Here we first formulate basic RL algorithms that will become the basis of our later models. In
110 RL, an agent learns to predict the expectation of future rewards associated with a given state, a
111 quantity termed as *value*¹¹. For simplicity, we will drop the dependency on time here, but note
112 that the basic results hold even if time is considered (Methods 1.1). Learning of value is driven
113 by RPEs (δ), the discrepancy between the actual and expected reward (r and V , respectively)
114 (Eq1). To improve the accuracy of the value prediction, RPEs are utilized to update the estimate
115 of V . This is done iteratively by adding a fraction (α) of δ (Eq2) where α defines the learning
116 rate.

117

$$118 \quad \delta = r - V \quad (1)$$

$$119 \quad V \leftarrow V + \alpha \cdot \delta \quad (2)$$

120

121 When the magnitude of reward r is fixed (i.e., deterministic environment), the value V learned
122 through this algorithm (Eq 1 and 2) converges on r and the RPE converges on zero. When the
123 magnitude of reward r varies stochastically trial-to-trial, the value at convergence fluctuates
124 around the expected value of the reward distribution (see Methods 1) (Fig. 1a) and the RPE
125 around zero.

126

127 **Risk-sensitive RL.** In the framework called risk-sensitive RL⁴⁰, learning rates are defined
128 separately for positive and negative RPEs (denoted by α^+ , α^-).

129

$$130 \quad V \leftarrow V + \alpha^+ \cdot \delta \text{ if } \delta > 0 \quad (3)$$

131 $V \leftarrow V + \alpha^- \cdot \delta$ if $\delta < 0$

132

133 In the presence of stochastic rewards, when the learning rates between positive and negative
134 RPEs are different, the value learned through this algorithm (Eq 1 and 3) does not converge on
135 the expected value, but instead on a value higher or lower than the expected value depending on
136 the relative amplitude of the learning rates α^+ , α^- . This algorithm, therefore, develops optimistic
137 or pessimistic value expectations, respectively. This learning algorithm is called “risk-sensitive”
138 because values of probabilistic (risky) rewards are biased compared to deterministic (certain)
139 rewards, and, therefore, the agent develops a preference between risky and certain rewards even
140 when the expected values are the same (Fig. 1b).

141

142 **Distributional RL.** The concept of asymmetric updates has been utilized in a novel RL
143 framework called distributional RL^{36,37,41}. This algorithm allows an agent to learn the entire
144 probability distribution of rewards, instead of the expected value which is typically the learning
145 target in traditional RL algorithms (Fig. 1c). In distributional RL, an agent is equipped with a set
146 of value predictors (V_i), where i corresponds to the index of the value predictor (or “value
147 neuron”). The value of the i -th neuron (V_i) is updated based on the learning rates (α_i^+ , α_i^-) and
148 the RPE (δ_i) for that neuron i :

149

150 $V_i \leftarrow V_i + \alpha_i^+ \cdot \delta_i$ if $\delta_i > 0$ (3)

151 $V_i \leftarrow V_i + \alpha_i^- \cdot \delta_i$ if $\delta_i < 0$

152

153 Similar to risk-sensitive RL, the learned value of each value predictor converges on estimates
154 larger or lower than the expected value, determined by the ratio between α_i^+ and α_i^- .

155 Mathematically, each V_i converges on the τ_i -th expectile of the distribution (Fig. 1c) where τ_i
156 (asymmetric scaling factor) is defined by:

157

158 Asymmetric scaling factor: $\tau_i = \frac{\alpha_i^+}{\alpha_i^- + \alpha_i^+}$ (4)

159

160 Expectiles are the solutions to asymmetric least squares minimization and generalize the mean of
161 a distribution (with the mean being the 0.5th expectile) as quantiles generalize the median (with
162 the median being the 0.5th quantile)⁴². Since a set of expectiles can define a distribution, the
163 diversity of τ_i across the population enables learning of the entire probability distribution.

164

165 **Problem.** In both risk-sensitive RL and distributional RL, unbalance in learning rate parameters
166 for positive and negative RPEs gives rise to optimistic and pessimistic biases in learned values.
167 Importantly, however, the underlying biological mechanism regulating learning rate parameters
168 (α^+ , α^-) and asymmetry thereof (τ) remains unclear.

169

170 In the following sections, we will discuss potential biological mechanisms that regulate
171 asymmetric learning rates (α^+ , α^-). We will first modify the above RL algorithms to incorporate
172 important neural circuit architectures in the brain. We will then propose two key biological
173 mechanisms that can give rise to asymmetric learning rates (called Model 1 and 2). We will then
174 show that our model can explain previous experimental data and psychiatric conditions.

175

176 **Incorporating biological features into RL models**

177 The above RL models provide algorithmic-level formulations, yet they do not recapitulate
178 fundamental characteristics of the neural circuits thought to perform RL in the brain⁴³⁻⁴⁶. We
179 next incorporate some of the important circuit and synaptic properties into the model.

180

181 In the brain, it is thought that dopamine neurons in the ventral tegmental area (VTA) broadcast
182 RPEs¹⁷ and modulate synaptic plasticity in dopamine-recipient areas^{33,47}. The striatum is the
183 major target of dopaminergic projections. It has been thought that spiny projection neurons
184 (SPNs) in the striatum represent values, and dopamine modulates plasticity of synapses on

185 SPNs^{33,34,47,48} (Fig. 2a). Under this framework, the value representations in SPNs are updated by
186 dopaminergic RPEs. In most RL models, each value predictor is typically updated by both
187 positive and negative RPEs. If the value is computed based on a weighted sum of some inputs
188 (i.e., using linear function approximation¹¹), the update rules described above (Eq 3 and 4) are
189 equivalent to performing a semi-gradient descent that minimizes RPEs¹¹ (see Methods).

190
191 The basic architectural assumptions of these RL models are, however, at odds with the RL
192 circuitry in the brain. Importantly, in the striatum, there are two major classes of dopamine-
193 recipient SPNs characterized based on the type of dopamine receptor that they express: D1- or
194 D2-type dopamine receptors (D1R and D2R)⁴⁸. SPNs expressing D1R and D2R constitute the so-
195 called direct and indirect pathways and exert opposing effects on downstream “output” neurons,
196 with each pathway promoting or opposing a certain output (e.g., movement).

197
198 In addition to the presence of direct and indirect pathways, there are two additional properties in
199 these opposing populations that need to be considered³². First, D1R and D2R have different
200 affinities to dopamine: high in D2R and low in D1R (EC₅₀ affinity constant is 1 μ M for D1R and
201 10 nM for D2R)^{49,50}. The dose-occupancy relationship of D1R and D2R are sigmoidal but they
202 are shifted with one another with respect to dopamine concentration (Fig. 2b). Importantly, at
203 normal dopamine levels (approx. 50-100nM)^{51,52}, D2Rs are mostly occupied while D1Rs are
204 mostly unoccupied (Fig. 2b). Although whether the affinities of D1R and D2R differ at the
205 molecular level has been questioned⁵³, a recent study showed that intracellular signaling through
206 protein kinase A (PKA) in D1- and D2-SPNs is triggered by a phasic increase and a decrease in
207 dopamine, respectively, in behaving animals³⁵. These results are consistent with (apparent)
208 difference in affinities of D1R and D2R assumed in previous studies⁴⁹, although the exact reason
209 for the difference remains to be clarified⁵³.

210
211 The second important property pertains to different learning rules in D1- and D2-SPNs which are
212 predicted from different affinities of the receptors. Consistent with the observed PKA signals in

213 these cells, recent studies have shown that glutamatergic inputs on D1-SPNs are potentiated by a
214 transient *increase* in dopamine, whereas those on D2-SPNs are potentiated by a transient
215 *decrease* in dopamine^{34,35} (Fig. 2c), supporting opposing plasticity rules between D1- and D2-
216 SPNs.

217

218 There have been previous efforts to incorporate in RL models the direct and indirect pathways
219 (also called “Go” and “NoGo” pathways, respectively) such as Opponent Actor Learning
220 (OpAL³⁰, OpAL*⁵⁴) and Actor learning Uncertainty (AU)³² models. These previous models
221 were developed as *Actor-Critic models*¹¹. Here, we will build on the AU model to focus on the
222 problem of value learning and extend it to support risk-sensitive RL and distributional RL. Our
223 model has two separate populations of value predictors corresponding to D1R- and D2R-SPNs,
224 that store the quantities P_i and N_i respectively (Eq 6, Fig. 2d). Mimicking dopamine's effect on
225 potentiation, P_i or N_i will increase their estimates if an RPE is positive or negative, respectively,
226 with the learning rates defined by α_i^+ , α_i^- (Eq. 6). Importantly, the value V_i can be obtained
227 simply by taking the difference between P_i and N_i . (Eq. 7).

228

229 D1R-SPN:

$$230 \quad P_i \leftarrow P_i + \alpha_i^+ \cdot |\delta_i| - \beta \cdot P_i \quad \text{if } \delta_i \geq 0 \quad (5)$$

$$231 \quad P_i \leftarrow P_i - \beta \cdot P_i \quad \text{if } \delta_i < 0$$

232 D2R-SPN:

$$233 \quad N_i \leftarrow N_i + \alpha_i^- \cdot |\delta_i| - \beta \cdot N_i \quad \text{if } \delta_i \leq 0$$

$$234 \quad N_i \leftarrow N_i - \beta \cdot N_i \quad \text{if } \delta_i > 0$$

$$235 \quad \text{Value: } V_i = P_i - N_i \quad (6)$$

236

237 where β is a decay parameter which represents synaptic decay in the absence of RPEs. This
238 model (Eq 6 and 7) preserves various essential properties of the previous RL models: (1)
239 learning in P and N can be combined to provide a simple update rule for value V , and (2) this
240 update rule approximates the gradient descent that minimizes RPEs (when $\beta = 0$, the update rule

241 is equivalent to the gradient descent). Importantly, with $\beta > 0$, we can show that these simple
242 learning rules guarantee convergence of value, without the need for additional mechanisms to
243 modulate the learning rates over iterations (Methods 1.3).

244

245 For instance, in a stochastic environment where there is a probability p of receiving a reward of a
246 fixed magnitude $r = 1$, the stochastic fixed point of the learned value V_i (i.e., convergence point)
247 will be defined by Eq 7.

248

$$249 \quad V_i = \frac{\frac{\tau_i \cdot p}{1-\tau_i} \cdot r}{\frac{\tau_i \cdot p}{1-\tau_i} + 1 + C} \cdot r, \text{ where } C = \frac{\beta}{(1-p) \cdot (1-\tau)} \quad (7)$$

250

251 Note that Eq.7 contains an additional term C which depends on β and this decay factor β is
252 important to stabilize the P_i and N_i estimates (avoid infinite increases) (Methods, 1.3.1, Extended
253 Data Fig. 1).

254

255 This formulation now provides a mechanistic model suitable for risk-sensitive RL (when there is
256 one value predictor) as well as distributional RL (when there are multiple value predictors),
257 which incorporate the neural circuit architecture and plasticity rules of D1R- and D2R-SNPs
258 found in the brain.

259

260 With this model at hand, we now discuss potential mechanisms that produce an asymmetry in
261 learning rates α_i^+ , α_i^- , which, in turn, causes biases in value predictions. In principle, learning
262 rate parameters can be a function of (1) the scaling of RPEs, i.e., the slope of dopamine
263 responses as a function of RPE (δ), and (2) the scaling of value updates, i.e., the efficacy of
264 dopamine-dependent synaptic plasticity at the level of SPNs. In the following, we discuss each
265 scenario, emphasizing the role of either tonic or phasic dopamine activity in each of these
266 mechanisms (Model 1 and 2, respectively). For simplicity, we will start with a model in which
267 α^+ , α^- are equal for all neurons within both P and N populations, equivalent to risk-sensitive

268 RL. We will then relax this assumption and introduce heterogeneity by allowing α_i^+, α_i^- to vary
269 across neurons, implementing a form of distributional RL.

270

271 **Model 1: The role of baseline dopamine in asymmetric learning**

272 As discussed above, D1R and D2R have different affinities to dopamine which leads to different
273 levels of receptors' occupancy at a given baseline dopamine level (Fig. 2b). Crucially, due to the
274 sigmoidal shape of the dose-occupancy curves, the slope of the curve changes with baseline
275 dopamine level, which means that a given dopamine transient leads to a different change in
276 receptor occupancy depending on the baseline dopamine level (Fig. 3a,b). That is, the receptors'
277 *sensitivity* changes with baseline dopamine (Fig. 3c). In addition, a key consequence of the
278 distinct receptors' affinities is that an increase and decrease in baseline dopamine will cause
279 opposite changes in the sensitivity of D1R and D2R. Specifically, an increase in dopamine will
280 decrease D1R sensitivity relative to D2R, whereas a decrease in dopamine will increase D2R
281 sensitivity relative to D1R (Fig. 3c,d).

282

283 Building on this insight in Model 1, we postulate that the learning rates for positive and negative
284 RPEs are a function of the D1R and D2R sensitivity, respectively. This is supported by previous
285 studies that have reported that the effect of dopamine transients of a given magnitude in SPNs'
286 plasticity can be modulated by the level of dopamine baseline³⁴. In addition, it has been reported
287 that the level of potentiation in SPNs^{33,34} or plasticity, which are related to intracellular signals³⁵,
288 scale with the magnitude of dopamine transients, keeping all else fixed. These observations can
289 be summarized with the following rule:

$$290 \quad LTP_{D1} \approx \alpha^+ \cdot |DA_{burst}|$$

$$291 \quad LTP_{D2} \approx \alpha^- \cdot |DA_{pause}|$$

292 Where α^+, α^- correspond to the receptors' sensitivities and depend on the dopamine baseline
293 level. This rule can be directly related to the update equations for the P and N populations in our
294 model:

$$295 \quad \Delta P = \alpha^+ \cdot |\delta| \dots \text{if } \delta \geq 0$$

296
$$\Delta N = \alpha^- \cdot |\delta| \dots \text{if } \delta < 0$$

297

298 It can be shown that this learning rule is in agreement with the normative solution for the
299 problem of value learning¹¹ (Methods 1.6).

300

301 In short, in Model 1, a shift in the baseline dopamine level causes asymmetries in scaling of the
302 value updates for positive versus negatives RPEs via the modulation of receptors' sensitivities,
303 which leads to value learning biases. This is a direct consequence of the dose occupancy
304 relationships of D1R and D2R (Fig. 3b-d).

305

306 **Model 2: Asymmetric scaling of phasic dopamine responses, inspired by distributional RL**

307 In Model 2, we postulate that the learning rates α^+ and α^- are a function of the scaling (i.e.,
308 'slope') of dopamine responses evoked by positive and negative RPEs, respectively:

309
$$DA_{burst} = \alpha^+ \cdot \delta \dots \text{if } \delta \geq 0$$

310
$$DA_{pause} = \alpha^- \cdot \delta \dots \text{if } \delta < 0$$

311 This is supported by a previous study on distributional RL that demonstrated that individual
312 dopamine neurons vary in terms of how the magnitude of reward responses is scaled as a
313 function of positive and negative RPEs (Fig. 3e)³⁶.

314

315 In the distributional RL framework, individual dopamine neurons vary in terms of their
316 asymmetric scaling factor τ_i and each of the multiple value predictors (V_i) converges on the τ_i -th
317 expectile of the reward distribution (Eq. 4). However, in most applications of distributional RL,
318 action selection is still based on the expected value of the reward distribution. Thus, the quantity
319 relevant to action selection can be described using the population level average $\tau_{population}$, and
320 biased value learning at the behavioral level could arise if $\tau_{population}$ is higher or lower than 0.5.
321 This can occur from a differential loss of optimistic or pessimistic dopamine neurons. Another
322 possibility is an overall upward or downward shift in the distribution of τ_i across the population
323 due to, for example, intrinsic factors modulating the gain of dopamine phasic responses.

324

325 Risk-sensitive RL can be thought of as a special case of distributional RL which has only one
326 value predictor. Here, the slope of the average dopamine evoked transient to positive and
327 negative RPEs, may correspond to the population level learning rates for positive and negative
328 RPEs (α^+ , α^-), respectively. If the asymmetric scaling factor τ is higher or lower than 0.5, value
329 learning will be biased (Fig. 3e).

330

331 **Testing for evidence of either model in experimental data**

332 *Tian and Uchida (2015).*

333 We next examined whether Model 1 or 2 can explain empirical data obtained in experimental
334 animals or humans. We first examined the data obtained in mice in our previous study⁵⁵. In this
335 study, the authors tested the effect of lesioning the habenula, a brain structure implicated in
336 depression⁵⁶⁻⁵⁸, on the activity of dopamine neurons and on reward-seeking behavior. Head-fixed
337 mice were trained in a Pavlovian conditioning task in which odor cues predicted reward with
338 different probabilities (10%, 50%, 90%). After performing habenula (n=5) or sham (n=7) lesions
339 (Fig. 4a), the spiking activity of VTA dopamine neurons was recorded while mice performed the
340 task.

341

342 After lesions, mice exhibited an elevated reward-seeking behavior (anticipatory licking) in
343 response to cues predictive of probabilistic rewards, consistent with an optimistic bias in reward
344 expectation (Fig. 4b, right). Importantly, anticipatory licking gradually increased over several
345 sessions after lesions, suggesting that the optimistic bias developed through learning (Fig. 4b,
346 left). To bring insight into the underlying cause of these biases, we fit two different RL models to
347 the anticipatory lick responses on a trial-by-trial basis (Extended Data Fig. 2), assuming a linear
348 relationship between value predictions and anticipatory licking. These models considered either a
349 change in the sensitivity to rewards (Extended Data Fig. 2b) or asymmetric learning rates
350 (Extended Data Fig. 2c). This analysis showed that the biases observed in the behavior could be
351 explained by asymmetric learning rates, but not by reward sensitivity because the reward

352 sensitivity was unchanged in the lesion group with respect to the control group (Extended Data
353 Fig. 2c).

354
355 Dopamine neurons' responses to reward-predictive cues reflect the increases in value expectation
356 predicted by the cue with respect to baseline. The overall magnitudes of cue-evoked responses
357 were not elevated in lesioned animals compared to control animals (Fig. 4d). However, the shape
358 of the response curve indicated an 'optimistic' bias: although in control animals, cue responses
359 scaled linearly with the expected value (i.e., reward probability), the response function of the
360 lesioned animals was convex. In other words, in control animals the response to the 50%-reward
361 cue was not significantly different from the quantity that results from the linear interpolation
362 between the responses to 10%- and 90%-reward cues. In lesioned animals, however, the response
363 to the 50%-reward cue was significantly greater than this quantity, which is indicative of an
364 optimistic bias in value predictions (Fig. 4d, see Methods 1.3.3 for analysis of value predictions
365 curve convexity). Such a change was observed at the level of the population average. Further
366 analysis using individual neurons showed that when calculating a single-cell level metric that
367 compares the 50%-reward cue to the same linear interpolation point, there was a broad
368 distribution in this metric below and above the interpolated point both in the control and lesion
369 groups (Fig. 4e-f). The distribution was, however, shifted in its mean in the lesion group (Fig.
370 4e). These analyses indicated that both anticipatory licking and dopamine cue responses have an
371 optimistic bias as characterized by an overvaluation of probabilistic rewards, without still
372 pointing to the underlying mechanism.

373
374 **Model 2 cannot explain the optimistic biases in behavior and cue-evoked dopamine**
375 **responses after Hb lesions**

376
377 In Model 2, an optimistic bias in reward expectation can arise if the average of the asymmetric
378 scaling factor at the population level ($\tau_{population}$) becomes greater than 0.5 (Fig. 5a,b).

379

380 To test this idea, we obtained the asymmetric scaling factors (τ_i) from dopamine neurons based
381 on their outcome responses: for each neuron, we constructed outcome response functions against
382 the magnitude of RPEs (Fig. 5c, Extended Data Fig. 3a,b). The response functions were obtained
383 based on (1) whether reward was delivered (positive RPEs) or not (negative RPEs), and on (2)
384 the magnitude of the reward expectation given by the reward probabilities predicted by each cue
385 (0.1, 0.5, 0.9) (Extended Data Fig. 3a,b). We then obtained the point at which the responses are
386 more likely to be below or above baseline (i.e., ‘zero-crossing points’)³⁶ (Extended Data Fig. 3c),
387 and computed α_i^+ and α_i^- as the slopes of the responses in the positive and negative domains
388 with respect to this zero-crossing point (Extended Data Fig. 3d), respectively. In both control and
389 lesioned animals, asymmetric scaling factors tiled a wide range between 0 and 1 and presented
390 other signatures consistent with distributional RL³⁶ (Extended Data Fig.4). Nonetheless, although
391 the variance of the distribution of asymmetric scaling factors was greater in lesioned animals, the
392 mean did not change, indicating a lack of bias between α_i^+ and α_i^- at the population level (Fig.
393 5d). This was also the case when the asymmetric scaling factor was derived directly from the
394 population average response (Fig. 5c). Thus, contrary to the conclusion in our previous study¹⁵,
395 these analyses indicated that changes in reward responses (and the resulting scaling factor τ) do
396 not explain the optimistic biases in behavior nor cue responses in lesioned animals (Fig. 5e,f).

397

398 **Model 1 can explain the optimistic biases in behavior and cue-evoked dopamine responses** 399 **by Hb lesion**

400 In addition to changes in the magnitude of dopamine RPEs, we observed that the baseline firing
401 rates of dopamine neurons were elevated in lesioned animals (Fig. 6a). According to Model 1, if
402 these changes are followed by an increase in the baseline dopamine levels in the striatum, this
403 should give rise to biased value learning ($\alpha^+ > \alpha^-$) and an optimistic bias in value expectation.
404 In this way, this change in baseline firing can explain optimistic biases observed in lesioned
405 animals. However, it remains unclear whether the observed change in baseline firing can result in
406 functionally relevant levels of changes in the receptor occupancies discussed above.

407

408 To quantitatively predict dopamine concentrations in the striatum and resulting receptor
409 occupancies of D1R and D2R, we used a biophysical model commonly used in the field⁵⁹
410 (Fig.6a). This model has the firing rate of dopamine neurons as its input, and considers diffusion
411 of dopamine, dopamine reuptake, and D2-autorreceptor-mediated inhibition of dopamine release
412 to predict the dopamine concentration in the striatum (Fig. 6b,e). In addition, it considers the
413 affinities of D1R and D2R to estimate their occupancy levels (Fig. 6c,f). After estimating these
414 two variables (dopamine concentration and receptor occupancy), we derived the receptor
415 sensitivities (Fig. 6g-h). The receptor sensitivities were quantified as the slope of the resultant
416 changes in receptor occupancy given the observed baseline and phasic responses of dopamine
417 neurons. We then trained Model 1 using the receptor sensitivities as learning rates (α^+ and α^-)
418 for both control and lesioned animals.

419
420 The biophysical model indeed supported that the observed change in dopamine neuron firing can
421 cause a significant increase in dopamine concentration (Fig. 6e) and in D1 and D2 receptor
422 occupancies at baseline (Fig. 6g). These changes are expected to cause a significant asymmetry
423 in receptor sensitivities favoring D1 receptors over D2 receptors (Fig. 6h-i).

424
425 These receptor sensitivities were directly used as the asymmetric learning rates in a temporal-
426 difference (TD) learning version of Model 1 (see Methods 1.3, 3.3). After training, the model
427 incorporating the predicted asymmetries in learning rates (α^+ , α^-) produced optimistic biases in
428 value predictions and in normalized cue responses, similar to those observed in lesioned animals
429 (Fig. 6k-l). The model simulating control animals developed no significant biases.

430
431 Additionally, the overall decrease in the magnitude of cue responses, observed in lesioned
432 animals, was reproduced in Model 1 using TD learning (Fig. 6k). This occurs because TD
433 learning calculates RPEs based on the change in values between before and after cue
434 presentation, and the “baseline” (pre-cue) reward expectation was also increased by optimistic
435 value learning (Fig. 6k). These results, together, indicate that Model 1 provides a parsimonious

436 account of the data: a change in baseline firing of dopamine neurons, rather than changes in
437 phasic responses, is the likely mechanism that led to optimistic biases in reward-seeking
438 behavior as well as cue-evoked dopamine responses in habenula lesioned animals.

439

440 **Model 1 and model 2 play complementary roles in the encoding of asymmetric learning**
441 **rates**

442 Although Model 2 did not explain the optimistic biases in the data in habenula-lesioned mice, the
443 distributional RL version of Model 2 explained other features of the data (Extended Data Fig. 3-
444 4). As mentioned, in both control and lesioned animals, asymmetric scaling factors tiled a wide
445 range between 0 and 1³⁶ (Extended Data Fig.4). Furthermore, cue-evoked responses of individual
446 neurons showed a wider distribution than what is expected by noise (Figure 4d). Finally, the core
447 prediction of distributional RL – a positive correlation between the asymmetric scaling factors of
448 the RPE responses of individual dopamine neurons and their zero-crossing points³⁶ – was also
449 present in controls and after Hb lesions. Together these results support that the basic features of
450 distributional RL are present in a way consistent with Model 2.

451

452 To complement this analysis, we tested whether Model 2 could have explained the signatures of
453 the data if asymmetric scaling factors (τ) derived from dopamine responses were indeed overall
454 biased (Extended Data Fig. 5). As expected from the model's fixed-point analysis (Methods 1.3),
455 if we imposed a shift in the mean of the distribution of asymmetric scaling factors (i.e.,
456 $\tau_{population} > 0.5$), the value predictors indeed exhibited optimistic biases (Extended Data Fig.
457 5e,f). However, the model did not reproduce the optimistic bias in cue-induced TD errors
458 observed in the data (Extended Data Fig.5g,h). This is due to an interaction of the biases in
459 prediction at “baseline” (pre-cue) and the cue, together with the optimistic asymmetry in the
460 scaling of the TD errors at cue themselves. Importantly, this was found in both versions of Model
461 2, distributional and risk-sensitive RL (Extended Data Fig. 5a-d and e-h). The difficulty of
462 explaining biased dopaminergic cue responses further makes the Model 2 an unlikely mechanism
463 to explain the optimistic biases in the data.

464

465 Altogether the data supports a model in which the mechanisms of Model 1 and 2 play
466 complementary roles in the encoding of asymmetric learning rates. The mechanism of Model 2
467 explains the variability in single neuron responses, consistent with the expectile code in
468 distributional RL. On the other hand, the mechanism of Model 1 at the population level,
469 generating asymmetries in learning rates and biases in value expectations, which might require
470 context-dependent regulation⁶⁰.

471

472 Taken together, the above results suggest that Model 1 and 2 coexist in the brain. This can be
473 formalized as follows:

474

$$475 \quad P_i \leftarrow P_i + \hat{\alpha}_i^+ \cdot \delta_i \dots \text{if } \delta_i \leq 0$$

$$476 \quad N_i \leftarrow N_i + \hat{\alpha}_i^- \cdot \delta_i \dots \text{if } \delta_i < 0$$

477 where:

$$478 \quad \hat{\alpha}_i^+ = \alpha_P \cdot \alpha_i^+$$

$$479 \quad \hat{\alpha}_i^- = \alpha_N \cdot \alpha_i^-$$

480

481 where α_i^+ and α_i^- correspond to the asymmetric scaling of dopamine RPEs at the single-cell level
482 (Model 2) and α_P and α_N correspond to the asymmetric scaling of synaptic plasticity in D1R and
483 D2R at the population level (Model 1).

484

485 **Linking asymmetric learning and baseline dopamine levels in healthy subjects**

486 *Cools et al., (2009)*²⁴.

487 There have been very few studies that examined the relationship between baseline dopamine
488 levels and asymmetry in learning from positive and negative outcomes. As a rare case for such
489 examinations, Cools et al.²⁴ provided intriguing data in humans. They compared the performance
490 in reversal learning and the quantity called ‘dopamine synthesis capacity’. Dopamine synthesis
491 capacity is estimated by injecting the positron emission tomography (PET) tracer

492 [¹⁸F]fluorometatyrosine (FMT) and is thought to be correlated with baseline dopamine levels^{61,62}.
493 This study found that higher dopamine synthesis capacity was correlated with better learning
494 from gains but not with learning from losses (Fig. 7b). As a result, in reversal learning, subjects
495 with higher dopamine synthesis capacity learned from gains than losses, reported as the ‘relative
496 reversal learning (RRL)’ index in their study (Fig. 7b). This result, thus, provides direct evidence
497 supporting our Model 1.

498
499 In addition, they found that dopamine synthesis capacity predicts the effectiveness of
500 bromocriptine (D2 partial agonist) in altering learning rate asymmetry: bromocriptine’s ability to
501 bias learning from gains over losses (i.e., positive change in RRL) was negatively correlated with
502 dopamine synthesis capacity (Fig. 7c). We found that this result can also be explained by Model
503 1. For this, we simulated the effects of bromocriptine with the biophysical model used above,
504 and derived the asymmetric learning rates from the slopes of the D2R occupancy (Fig. 7d,
505 Extended Data Fig. 6a,b) or activation curves (Fig. 7d, Extended Data Fig. 6c,d). The RRL
506 parameter reported by Cools et al. corresponds to the asymmetric scaling factor τ , and is
507 equivalent to $(2\tau - 1)$ (as described in the Methods 4.1). We then computed what would be the
508 change in this parameter $\Delta(2\tau - 1)$ induced by bromocriptine (Fig. 7e-f, Extended Data Fig. 6e-
509 1).

510
511 This analysis revealed that by considering the asymmetries in learning rates induced by changes
512 in the baseline occupancy of the receptors, our model can capture their results in a qualitative
513 manner. Intuitively, the less dopamine there is at baseline, the lower the occupancy of D2R at
514 placebo conditions. This leads to a larger increase in D2R occupancy induced by D2 agonist in
515 low dopamine baseline conditions (Fig. 7d, Extended Data Fig. 6a) and, thus, a larger increase in
516 asymmetry in learning from gains over losses, if D1R occupancy is kept fixed. These effects still
517 hold even if we consider, in addition to bromocriptine’s effects in postsynaptic receptors (D2
518 long or D2l), its effect on inhibition of dopamine release via presynaptic (D2 short or D2s)
519 autoreceptors^{63,64} (Fig. 7d, Extended Data Fig. 6b). This can be simulated as a decrease in

520 dopamine level, which leads to a shift in the occupancy curves to the right. Finally, we can
521 consider effect of the *partial* agonism of the drug, that leads to a lower activation level of
522 receptors even if the occupancy is maximal (Fig. 7d, Extended Data Fig. 6c-d). Even after
523 considering this last factor, the results remain qualitatively the same as those found in the
524 original study. These results were robust to a relatively wide range of values in the simulation's
525 parameters (Extended Data Fig. 7, 8).

526

527 **Linking psychiatric conditions to baseline dopamine levels**

528 *Timmer et al., 2018*

529 Various psychiatric disorders are characterized by abnormal future predictions or mood. Our
530 Model 1 raise the possibility that an overall decrease in baseline dopamine level in the striatum
531 would enhance learning from negative outcomes over learning from positive outcomes leading to
532 persistent pessimistic future value expectations, a hallmark of depressive-like symptoms (Fig.
533 3a,b). A piece of evidence supporting this in the human literature is the greater learning rates for
534 losses over gains in patients with Parkinson's disease (PD)^{22,25}, its comorbidity with
535 depression^{25,65} that can precede the PD diagnosis⁶⁵⁻⁶⁷, and the reports of decreased dopamine
536 transporter binding in the ventral striatum in depressed PD patients compared to non-depressed
537 PD patients^{68,69}.

538

539 In addition, the progression of dopaminergic axonal loss in PD is topographically unbalanced:
540 the axonal loss is more prominent in the dorsal striatal regions⁷⁰ than in the ventral ones. This
541 leads to uneven dopamine baseline levels across the striatum that would interact with the global
542 increases in dopamine induced by dopaminergic medications in PD patients. We hypothesize that
543 a behavioral readout of the degree of this unevenness might be the presence or absence of
544 depression as a comorbidity: *patients with depression might have lower dopamine levels in the*
545 *ventral striatum*. Thus, if indeed baseline dopamine levels are correlated with depression, this
546 comorbidity could be predictive of the effects of PD medication.

547

548 We examined a previous study that provided evidence for this hypothesis ⁷¹. Here, PD patients
549 with and without depression history were tested in a gambling task, under presence or absence of
550 medication ('ON' and 'OFF' medication states). The authors fitted a 'loss aversion' parameter to
551 the behavioral performance, which is equivalent to $1 - \tau$ in our model, under some assumptions
552 (Methods). Their results were consistent with our model predictions. In the OFF-medication
553 state, there was a (near-significant) main effect of depression group (with or without depression)
554 on the learning rate asymmetry: patients with a depression history tended to be more loss averse
555 than nondepressed patients ($P = 0.052$). This is consistent with a decrease of dopamine levels in
556 the ventral striatum and thus a regime of $\alpha^+ < \alpha^-$ in value learning. Importantly, in the ON-
557 medication state, the medication effects on the asymmetry in learning rates were predicted by the
558 degree of severity of depression: patients with larger depression scores exhibited greater drug-
559 induced decreases in loss aversion (Fig. 7g), which would correspond to an increase in $\tau =$
560 $\frac{\alpha^+}{\alpha^+ + \alpha^-}$ in our model This is consistent with our Model 1: higher degrees of depression might be
561 correlated with lower levels of baseline dopamine, making the D1R sensitivity more susceptible
562 to an artificial increase in baseline dopamine with L-DOPA medication (Fig. 7h; further details
563 discussed in Methods).

564

565

566 **Discussion**

567

568 A hallmark of various psychiatric disorders is overly optimistic or pessimistic predictions about
569 the future. Using RL models, we sought to identify potential biological mechanisms that give rise
570 to biased value predictions, with a particular focus on the roles of phasic versus tonic dopamine.
571 Our results demonstrate that variations in tonic dopamine levels can modulate the efficacy of
572 synaptic plasticity induced by positive versus negative RPEs, thereby resulting in biased value
573 learning (Model 1). This effect arises due to sigmoidal shapes of the dose-occupancy curves and
574 different affinities of dopamine receptors (D1R and D2R); alterations in the tonic dopamine level

575 result in changes in the slope of the dose-occupancy curve (and thus, sensitivity) of dopamine
576 receptors at the baseline dopamine concentration. We show that this mechanism offers a simple
577 explanation for how changes in tonic dopamine levels can result in biased value learning in a few
578 examples of value learning in mice and humans. Additionally, we show that this mechanism may
579 underlie symptoms of various psychiatric and neurological disorders. Although altered phasic
580 dopamine responses could have been a natural suspect as a candidate mechanism for biased
581 value learning^{37,38}, our study provides a novel mechanism; the interaction between tonic and
582 phasic dopamine can give rise to biased value learning, even when phasic dopamine responses
583 remain relatively unchanged.

584

585 **The impact of properties of dopamine receptors on reinforcement learning (RL)**

586 Our results highlight the importance of considering properties of dopamine receptors and neural
587 circuit architecture (i.e., direct and indirect pathways) in RL models. Based on different affinities
588 of dopamine D1 and D2 receptors, it has been proposed that D1- and D2-SPNs play predominant
589 roles in learning from positive and negative dopamine responses^{32,72–75}. In support of this idea,
590 recent experiments have demonstrated that PKA signaling in D1- and D2-SPNs is primarily
591 driven by a phasic increase and decrease of dopamine, respectively³⁵. Furthermore, LTP-like
592 changes in D1- and D2-SPNs are triggered by a phasic increase and decrease of dopamine,
593 respectively^{33,34}. These recent pieces of evidence suggest that these plasticity rules are a basic
594 principle of the RL circuitry in the brain. Here we explored the properties of this RL model and
595 found the impact of the shape (slope) of receptor occupancy curves and showed that the tonic
596 dopamine levels can modulate the relative efficacy of learning from positive versus negative
597 RPEs.

598

599 One assumption in our model is that after a change in the tonic dopamine level, intracellular
600 signaling reaches a steady inactive state, and it is the *change* in receptor occupancy that matters
601 for inducing synaptic plasticity, rather than the *absolute* level of receptor occupancy reached
602 during phasic dopamine responses. We note that absolute level might also contribute, yet it is

603 expected that an increase or decrease in absolute occupancy levels will cause effects in the same
604 direction as the effects of relative change that we explored in this study.

605
606 Additionally, our model, which incorporates the new plasticity rules, the opponent circuit
607 architecture and properties of D1/D2 dopamine receptors, provides insights into the basic design
608 principle of the brain's RL circuit. It should be noted that the dose occupancy curves were
609 plotted as a function of the logarithm of dopamine concentration, which makes the occupancy
610 curves into sigmoidal shapes (Fig. 3, Extended Data Fig. 9). This logarithmic scaling is
611 important in two ways. First, considering two sigmoidal curves for D1R and D2R together, the
612 curves are approximately *symmetric* around the normal baseline dopamine level (Fig. 3a,
613 Normal). Second, logarithmic scaling means that a fold-change in dopamine concentration will
614 lead to the same leftward or rightward shift in these plots. It has long been argued that signaling
615 of RPEs by dopamine neurons is curtailed by the fact that dopamine neurons have relatively low
616 firing rates (2-8 spikes per second), and inhibitory responses of dopamine neurons tend to be
617 smaller than excitatory responses^{76,77}. Importantly, if we consider logarithmic scaling of
618 dopamine concentration, the problem of this asymmetry is substantially mitigated (Extended
619 Data Fig. 10). For example, with the baseline firing of 6 spikes per second, a phasic increase to
620 18 spikes per second and a phasic decrease to 2 spikes per second will cause the identical *fold-*
621 *changes* in spiking (i.e., 3-fold changes in both directions), which would lead to a similar *fold-*
622 *changes* in dopamine levels (Extended Data Fig. 11) and similar *percent* increase and decrease in
623 receptor occupancy in D1R and D2R, respectively (Fig. 3a). Consequently, the system achieves
624 symmetry in its response to positive and negative dopamine responses of observed magnitudes.

625
626 This may help understand *why* the basal ganglia circuit employs the opponent circuit architecture
627 in the first place. In the model used in the present study, the value is encoded as the difference
628 between the activity of D1- and D2-SPNs ($V = P - N$)³². We propose that this opponent circuit
629 architecture, together with the logarithmic scaling of dopamine concentration, allows the system
630 to effectively learn and encode both positive and negative values, which are contributed by the

631 increase of firing in D1- and D2-SPNs, respectively. This would allow to expand the dynamic
632 range of value coding, without requiring high baseline firing rates. Thus, at the normal dopamine
633 baseline, learning from positive and negative dopamine responses is well balanced. When the
634 tonic dopamine level deviates from the normal level, however, then the symmetry is broken and
635 value learning becomes biased, as explored in the present study.

636

637 **The role of tonic dopamine levels in psychiatric disorders**

638 As mentioned above, our modeling results provide an account for biased value predictions
639 observed in various psychiatric and neurological conditions. For one, our model provides a link
640 between findings in depressive-like states in animal models and the value learning biases
641 exhibited by humans.

642

643 In a rodent model of depression, it has been reported that spontaneous activity of dopamine
644 neurons is decreased⁷⁸ (but see^{79,80}). In addition, decreased spontaneous firing of dopamine
645 neurons has been observed as a result of chronic pain-induced adaptations that correlate with
646 anhedonia-like behavior⁸¹. Furthermore, maternal deprivation, which increases susceptibility to
647 anhedonia, led to an upregulation of D2R expression in the VTA⁸², which is expected to decrease
648 the excitability of dopamine neurons via its autoreceptor function. Finally, chronic
649 administration of corticosteroids, a method to mimic anxiety and anhedonia-like states, results in
650 an increase in somatodendritic dopamine concentration which then decreases dopamine
651 excitability via D2R hyper-activation⁸³. These results of decreased dopamine excitability
652 correlated with anhedonia-like states are consistent with findings of increased burst firing of
653 lateral habenula (LHb) neurons⁵⁶ and potentiation of glutamatergic inputs onto the habenula⁵⁷ in
654 depression models. This is further supported by reports that depressive-like behavioral
655 phenotypes can be ameliorated by optogenetic activation of dopamine neurons⁸⁴ and the anti-
656 depressant effects of ketamine might be mediated by the inhibition of bursting in the LHb⁵⁸

657

658 The mechanism by which a broad change in dopamine excitability could lead to depressive-like
659 states remains to be revealed. Just by assuming that a decrease in spontaneous firing leads to a
660 decrease in baseline dopamine level in the striatum, our model readily predicts that learning from
661 negative outcomes will be emphasized over learning from positive outcomes (Fig. 3a,b), as has
662 been reported in some studies of patients with major depressive disorder (MDD)¹. In addition,
663 RL agents learning in these conditions exhibit enhanced risk-averse behavior, pessimistic
664 outcome expectations, and increased sensitivity to losses compared to gains, all of which are
665 signatures of depressive-like conditions^{1,5,21,85,86}. This contrasts with findings of increased
666 dopamine synthesis capacity in pathological gambling patients⁸⁷, who show the opposite
667 behavioral signatures³.

668
669 An additional line of research relevant to our proposal is PD patients and pathological gambling
670 as a comorbidity. Previous work has emphasized the interaction between the degree of
671 dopaminergic loss and the effects of PD medications⁸⁸⁻⁹⁰, which can sometimes result in the
672 development of addictive disorders such as pathological gambling. As mentioned, the loss of
673 dopaminergic axons in PD patients has been reported to happen predominantly in the dorsal
674 regions of the striatum⁷⁰. Thus, at the onset of the motor impairment symptoms, which is when
675 L-DOPA medication tends to be prescribed, dopamine level is expected to be low in the dorsal
676 striatum while it might be relatively intact in the ventral striatum. This can lead to ‘overdose’ of
677 dopamine by medication: while L-DOPA might take dopamine levels in the dorsal striatum back
678 to its original set-point, it might cause an ‘overdose’ in the ventral striatum^{89,91}. Our model
679 predicts that this overdose would lead to decreases in D2R sensitivity relative to D1R. Assuming
680 that the ventral striatal regions have a dominant role in value learning, this would result in
681 excessive optimistic expectations and risk seeking, two key behavioral features of pathological
682 gambling and addictive disorders. We provided indirect evidence for this hypothesis; future work
683 should directly test these predictions.

684

685 It should be noted that we did not consider changes in dopamine receptors density, which have
686 also been related to value learning biases⁹² and psychiatric conditions⁹³. Future studies should
687 explore the influence of this additional factor in the encoding of asymmetric learning rates (i.e.,
688 $(\hat{\alpha}_i^+, \hat{\alpha}_i^-)$).

689

690 **Tonic dopamine as a modulator of ‘mood’**

691 Mood refers to a person’s emotional state as it relates to their overall sense of well-being.

692 Although the exact neural substrate of mood remains unknown, recent studies have indicated that
693 mood reflects not the absolute goodness of outcomes but rather on the discrepancy between
694 actual and expected outcomes in recent history^{13,14}. That is, mood depends on the cumulative
695 sum of RPEs that occurred recently¹³. It has also been proposed that mood, in turn, affects the
696 way we perceive and learn from positive and negative outcomes (RPEs)¹³.

697

698 Our model provides a unified mechanism for these two aspects of mood; both subjective feeling
699 of mood and biased learning from positive versus negative outcomes can arise from changes in
700 baseline dopamine levels which can be modulated by recent history of phasic dopamine
701 responses. It was proposed that this history dependent modulation of learning is an adaptive
702 mechanism that allows organisms to adapt quickly to slow changes in environments based on the
703 “momentum” of whether the situation is changing in a better or worse direction on a slow
704 timescale (e.g. seasonal change)^{13,14}. The models presented in the present study may provide
705 mechanistic insights into such mood-dependent modulation of learning and perception.

706

707 **Neural circuits for distributional reinforcement learning (RL)**

708 We examined the possibility that optimistic biases in reward seeking behavior and dopamine cue
709 responses observed in habenula-lesioned mice can be explained by Model 2, either based on risk-
710 sensitive RL (the average response) or distributional RL (responses of a diverse set of individual
711 dopamine neurons). We did not find evidence supporting this possibility. However, the present
712 study makes two important contributions with respect to distributional RL. First, we can show

713 that our model, which incorporated direct and indirect pathway architecture, can support
714 distributional RL (Extended Data Fig. 12, 13). It would be interesting to examine what additional
715 features and functions could be gained by having this opponent architecture. Second, we largely
716 replicated the previous results³⁶ using an independent data set. That is, the signatures of
717 distributional RL were present in this data set (Extended Data Fig. 3-4), and dopamine cue-
718 evoked responses did show an optimistic bias. This provides further evidence for a distributional
719 code in dopamine neurons, and shows that there is an overall elevated distributional
720 representation in dopamine cue responses in habenula lesioned animals.

721

722 **Concluding remarks**

723 Taken together, our biologically inspired RL model provides a foundation to link findings in the
724 brain and formal models of RL. Our work highlights a causal impact of baseline dopamine on
725 biasing future value predictions, which may underlie mood and some abnormalities observed in
726 psychiatric patients and could be used to regulate risk sensitive behavior.

727

728

729 **Methods**

730 **1. Reinforcement learning model**

731 Here we provide formal definitions and the framework of reinforcement learning used in this study. We
732 have focused our model formulations to the problem of *prediction*, in which an agent learns to predict the
733 value function¹¹. The problem of *control* (the problem of how an agent selects and executes actions) is not
734 considered. In RL, an agent's objective is to maximize the total cumulative rewards. It does so by learning
735 the value associated with each state in an environment. For now, we will develop the model dropping the
736 dependency on time within each episode. Here, the target to learn is the value function as defined by

$$737 \quad V(s_i) := \mathbf{E}[r^{(n)} | s^{(n)} = s_i]$$

738 Where $r^{(n)}$ is the reward experienced in the episode n (i.e., trial) of visiting state s_i . Learning of $V(s)$ is
739 driven by reward prediction errors (RPEs, δ), the discrepancy between the actual and expected reward:

$$740 \quad \delta^{(n)} = r^{(n)} - V(s_i)$$

741 The value is updated for the experienced state according to:

$$742 \quad V^{(n+1)}(s_i) \leftarrow V^{(n)}(s_i) + \alpha \cdot \delta^{(n)}$$

743 This is also known as the Rescorla-Wagner (RW) delta rule⁹⁴. The reward in each trial is sampled from a
744 reward distribution specific to a given state: $r^{(n)} \sim R(s_i)$. With the learning rule above, the value
745 converges on the expected value of this reward distribution. This can be shown with a stochastic fixed-
746 point approach; the convergence point is derived by obtaining the value of $V(s_i)$ at which the change in
747 $V(s_i)$ from trial n to trial $(n + 1)$ is expected to be zero (i.e., is zero on average):

$$748 \quad \mathbf{E}[V^{(n+1)}(s_i) - V^{(n)}(s_i)] = 0$$

$$749 \quad \mathbf{E}[\alpha \cdot \delta^{(n)}] = 0$$

$$750 \quad \mathbf{E}[\alpha \cdot (r^{(n)} - V(s_i))] = 0$$

$$751 \quad \alpha \cdot \mathbf{E}[r] - \alpha \cdot \mathbf{E}[V(s_i)] = 0$$

$$752 \quad \mathbf{E}[V(s_i)] = \mathbf{E}[r_t]$$

$$753 \quad V^*(s_i) = \mathbf{E}[r_t]$$

754 Where $V^*(s_i)$ is the stochastic fixed-point: the value around which $V(s_i)$ is expected to fluctuate after
755 learning and corresponds to the learning target above.

756 **1.1. Temporal difference learning**

757 Now we will consider time and extend the models to the temporal difference (TD) learning framework¹¹.
758 Dopamine responses have been shown to present key signatures of TD errors⁹⁵. Therefore, TD learning
759 models allow us to directly link the model variables to dopamine neural responses.

760 We can derive TD learning by defining a different environmental structure and learning objective. We
761 start by considering arbitrary states (s_t), which transition at each time step following a Markov process,
762 and at each time step the agent samples a random reward from a probability distribution $r_t \sim R(s_t)$.

763 The learning objective is now the value of a given state $V(s_t)$ defined as the *expected cumulative sum of*
764 *all future rewards* starting from state s . Rewards are discounted by a constant discounting factor (γ , with
765 $0 \leq \gamma \leq 1$) each time step. The expectation is taken over stochastic state transitions and sampled rewards:

$$766 \quad V(s_t) := \mathbf{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} \dots | s_t = S]$$

767 Where s_t is the state at time t , r_t is the reward sampled at time t and $V(s_t)$ is the value of the state s_t .
768 Since the environment and transitions are assumed to follow a Markov process, the equation above can be
769 rewritten in a recursive manner. This is known as the Bellman equation¹¹ :

$$770 \quad V(s_t) := \mathbf{E}[r_t + \gamma \cdot V(s_{t+1}) | s_t = S]$$

771 The agent approximates the true value $V(s_t)$ with a learned estimate $\hat{V}(s_t)$. With this approximation,
772 before learning converges, the estimates for the left- and right-hand sides are not equal. Thus, after
773 sampling a reward $r_t \sim R(s_t)$ from the environment, the difference between the two terms in the Bellman
774 equation represents the error in value prediction, called the temporal difference reward prediction error
775 (TD RPE, δ below),

$$776 \quad \delta_t = r_t + \gamma \cdot \hat{V}(s_{t+1}) - \hat{V}(s_t)$$

777 With α as the learning rate, the updates for the value estimates are:

$$778 \quad \hat{V}(s_t) \leftarrow \hat{V}(s_t) + \alpha \cdot \delta_t$$

779 With this definition, the TD RPE contains the difference between the estimated value of states evaluated
780 at consecutive time points. If we fix the discounting factor to be $\gamma = 1$, then $\gamma \cdot \hat{V}(s_{t+1}) - \hat{V}(s_t)$ is the

781 temporal derivative of the value function. As a result of this property, unexpected increases and decreases
782 in value result in positive and negative transient changes in TD RPE, respectively⁹⁵.

783 If dopamine responses encode TD RPEs, then cue-evoked responses can be formulated as:

$$784 \quad \delta_{cue} = \gamma \cdot \hat{V}(s_{cue}) - \hat{V}(s_b)$$

785 Where δ_{cue} is the TD RPE induced by the cue, $\hat{V}(s_b)$ is the value prediction at baseline and $\hat{V}(s_{cue})$ is
786 the value prediction elicited by the cue (which reflects the expected value predicted by each trial type). As
787 the $\hat{V}(s_b)$ is the same across all trial types and represents the average value predictions across them, then
788 δ_{cue} is dominated by the expected value of each trial type. This is a useful feature that we used in our
789 simulations for the habenula lesion experiment.

790 1.2. Distributional TD learning

791 In Results, we used a distributional TD learning model to test whether the subtle changes in the
792 distribution of asymmetric scaling factors observed after lesions could lead to the observed changes in cue
793 responses after learning.

794 In distributional TD learning, our learning objective is the entire distribution over cumulative discounted
795 future rewards, instead of the value defined above^{36,37,39}. We will call this the *return distribution*, $Z(s_t)$.
796 We can thus write an analogue of the Bellman equation, the ‘distributional Bellman equation’:

$$797 \quad Z(s_t) := R(s_t) + \gamma \cdot Z(s_{t+1})$$

798 The target to learn in distributional TD is now $V_i(s_t)$ that minimizes for the expectile regression loss:

$$799 \quad V_i(s_t) := \operatorname{argmin}_v E[(Z(s_t) - v)^2 \cdot (\tau_i - \mathbf{1}_{(Z(s_t) - v) < 0})]$$

800 Where $Z(s_t)$ is a random variable, representing the return distribution, and $\mathbf{1}_f$ is the indicator functions
801 that is equal to 1 if the condition in the subscript $\{f := (Z(s_t) - v) < 0\}$ is met, and 0 otherwise.

802 Minimizing the expectile regression loss makes $V_i(s_t)$ to converge on the τ_i^{th} expectile of the return
803 distribution³⁹.

804 The target is learned by taking samples from the estimated return distribution³⁹ $\tilde{z}(s_{t+1}) \sim Z(s_{t+1})$ and
805 from the reward distribution $r_t \sim R(s_t)$, to compute the TD error:

$$806 \quad \delta_{i,t} := r_t + \gamma \cdot \tilde{z}(s_{t+1}) - V_i(s_t)$$

807 Note that $\tilde{z}(s_{t+1})$ is random so the TD error is also random, and $\delta_{i,t} \neq r_t + \gamma \cdot V_i(s_{t+1}) - V_i(s_t)$. For
 808 more information regarding the sampling method employed in the simulations see Methods Section 3.3.
 809 In addition, the updates are performed with different learning rates (α_i^+ , α_i^-) for positive and negative δ_i .
 810 This asymmetry in the weighting of the errors used to update $V_i(s_t)$ is essential to minimize the expectile
 811 regression loss.

$$812 \quad \hat{\delta}_{i,t} = \alpha_i^+ \cdot \delta_{i,t} \dots \text{ if } \delta_{i,t} > 0$$

$$813 \quad \hat{\delta}_{i,t} = \alpha_i^- \cdot \delta_{i,t} \dots \text{ if } \delta_{i,t} < 0$$

814 The reliance on a single sample for $\tilde{z}(s_{t+1})$ suffers from high variance. Therefore, for performing the
 815 updates we average across a set of M updates, each depending on a single sample $\delta_{i,t}$.

$$816 \quad \mathbf{E}[\Delta V_i(s_t)] = \frac{1}{M} \sum_j^M \alpha_i^- \cdot \delta_{i,j} \cdot \mathbf{1}_{\delta_{i,j} < 0} + \alpha_i^+ \cdot \delta_{i,j} \cdot \mathbf{1}_{\delta_{i,j} > 0}$$

$$817 \quad V_i(s_t) \leftarrow V_i(s_t) + \mathbf{E}[\Delta V_i(s_t)]$$

818 This learning rule will asymptotically converge to the τ_i -th expectile of the return distribution³⁹.

819 1.3. TD learning with D1 and D2 populations

820 It is straightforward to extend the TD learning algorithm to have separate populations for D1 and D2
 821 SPNs³². We employed this model to derive dopamine cue responses with Model 1 (Fig. 6i). In this model,
 822 the same computation of TD RPE of standard TD learning is still used. Yet, this model differs in the
 823 updates and computation of $\hat{V}(s_t)$.

824 As mentioned previously, the updates in the P_i and N_i populations happen exclusively with positive or
 825 negative TD RPEs, respectively:

$$826 \quad P(s_t) \leftarrow P(s_t) + \alpha^+ \cdot |\delta_t| - \beta \cdot P(s_t) \dots \text{ if } \delta_t > 0$$

$$827 \quad N(s_t) \leftarrow N(s_t) + \alpha^- \cdot |\delta_t| - \beta \cdot N(s_t) \dots \text{ if } \delta_t < 0$$

828 Where α^+ and α^- are the learning rates for the P and N populations, that we postulate is modulated by
 829 baseline dopamine levels. The variable $\beta \in (0,1)$ is the decay factor, which we keep constant throughout
 830 the simulations and serves to stabilize $P(s_t)$, $N(s_t)$.

831 The computation of value estimate $\hat{V}(s_t)$ is given by:

832
$$\hat{V}(s_t) = P(s_t) - N(s_t)$$

833 **1.3.1. Convergence of risk sensitive TD learning**

834 We now discuss the convergence of the proposed TD learning algorithm with D1 and D2 populations.
 835 This analysis builds on the work in risk-sensitive reinforcement learning⁴⁰ and the already established
 836 results of convergence for stochastic iterative algorithms (e.g., TD learning) (Bertsekas & Tsitsiklis,
 837 1996⁹⁶, Proposition 4.4, p. 156) .

838 **Theorem:** The results by Bertsekas & Tsitsiklis (1996)⁹⁶ establish that, given a sequence $r_t \in$
 839 \mathbb{R}^m generated by the iterative algorithm:

840
$$a_{n+1}(s) = (1 - \sigma_n(s))a_n(s) + \sigma_n(s)((Ha_n)(s) + \omega_n(s)) \quad \forall s \in 1, \dots, m \quad \mathbf{Eq. I}$$

841 The variable a_n converges to the unique solution a^* of the equation: $Ha^* = a^*$ with probability = 1,
 842 assuming the following conditions are fulfilled:

843 1. The step sizes $\sigma_i(i)$ are non-negative and satisfy:

844
$$\sum_{n=0}^{\infty} \sigma_n(s) = \infty \quad \forall s \in 1, \dots, m$$

845
$$\sum_{n=0}^{\infty} \sigma_n(s)^2 < \infty \quad \forall s \in 1, \dots, m$$

846 2. The noise term $\omega_n(s)$ satisfies:

847 - $E[\omega_n(s)|\mathcal{F}_n] = 0 \quad \forall s, n$, where \mathcal{F}_n denotes the history of the process up to and including time
 848 step n

849 - Given any norm $\|\cdot\|$ on \mathbb{R}^m there exist constants A and B such that: $E[\omega_n^2(s)|\mathcal{F}_n] \leq A +$
 850 $B\|r_n\|^2 \quad \forall s, n$

851 3. The mapping H is a *maximum norm contraction* (see below for definition)

852 To prove convergence, we will first discuss the case of risk-sensitive TD learning following⁴⁰ and then
 853 discuss TD learning with D1 and D2 populations.

854 We define the risk sensitive TD-learning rule as:

855
$$\hat{V}_n(s) \leftarrow \hat{V}_{n-1}(s) + \sigma \cdot \mathcal{X}^\tau(\delta_{s_{n-1}, s_n})$$

856 Where:

857 $\delta_{s_{n-1}, s_n} = r_{n-1, n} + \gamma \cdot \hat{V}_{n-1}(s_n) - \hat{V}_n(s_{n-1})$

858 - The step index is $n \in 0, \dots, \infty$

859 - The step size σ is kept constant across iterations.

860 - For simplicity in calculations we follow ⁴⁰ and make use of the operator \mathcal{X}^K with $K \in (-1, 1)$

861
$$\mathcal{X}^K(x) = (1 - K) \cdot x \quad \dots \text{if } x > 0$$

862
$$\mathcal{X}^K(x) = (1 + K) \cdot x \quad \dots \text{if } x \leq 0$$

863 It is simple to show that the asymmetric scaling factor used in this paper is a scaled version of
864 the operator. That is: $\tau = 0.5(1 - K)$ and $1 - \tau = 0.5(1 + K)$.

865 - In addition, as in ⁴⁰, given that the function $\mathcal{X}^\tau(x)$ is piece-wise differentiable we can apply
866 the mean value theorem to show that for each pair of numbers (a, b) there exists a $\varepsilon_{a,b,K} \in$
867 $[1 - |K|, 1 + |K|]$, such that: $\varepsilon_{a,b,K} = \frac{\mathcal{X}^\tau(a) - \mathcal{X}^\tau(b)}{a - b}$. This relationship will become useful in
868 the future.

869 We will re-format the update rule to better match the iterative algorithm above:

870 Adding and subtracting $\sigma \cdot \hat{V}_{n-1}(s) / \alpha$

871
$$\hat{V}_n(s) \leftarrow (1 - \sigma/\alpha) \hat{V}_{n-1}(s) + \sigma/\alpha \left(\alpha \cdot \mathcal{X}^\tau(\delta_{n-1}) + \hat{V}_{n-1}(s) \right)$$

872 Defining an operator that will become useful:

873
$$\mathcal{T}_{\alpha K}[V](s) := V(s) + \alpha \cdot \sum_{i,j \in \mathcal{S}} p_{ij} \cdot \mathcal{X}^K \cdot \delta_{ij}$$

874 Defining the noise term as:

875
$$\omega_{n-1}(s) = \hat{V}_{n-1}(s) + \alpha \cdot \mathcal{X}^K(\delta_{s_{n-1}, s_n}) - \mathcal{T}_{\alpha K}[\hat{V}_{n-1}](s)$$

876 Then our update rule above becomes:

877
$$\hat{V}_n(s) \leftarrow (1 - \sigma/\alpha) \hat{V}_{n-1}(s) + \sigma/\alpha \left(\mathcal{T}_{\alpha K}[\hat{V}_{n-1}](s) + \omega_{n-1}(s) \right)$$

878 The formulation above can be directly compared to the one of stochastic iterative algorithm theorem

879 (Eq.I), and now we can check whether the conditions for convergence are met.

- 880 1. The conditions for the learning rate, are a direct consequence of our choice of the parameter
 881 which is a constant in our model and $0 < \alpha < 1$.
- 882 2. It has been shown that showed that the conditions for the noise term $\omega_{n-1}(s)$ as formulated
 883 above are satisfied⁴⁰.
- 884 3. Finally, the operator $\mathcal{T}_{\alpha\tau}[V](s)$ is a contraction mapping as also shown in Bersekas (1996)⁴⁰.

885 Therefore, the variable V_n converges to the unique solution V^* for which:

$$886 \quad V^* = \mathcal{T}_{\alpha K}[V^*] = V^* + \alpha \cdot \sum_{i,j \in \mathcal{S}} p_{ij} \cdot \mathcal{X}^K \cdot \delta_{ij}$$

887 We elaborate now on the proof for the contraction mapping of the operator $\mathcal{T}_{\alpha\tau}[V](s)$, as this will be
 888 useful for the proof of the D1 D2 TD learning model.

889 **Definition of contraction mapping.** Let (X, d) be a metric space (a set X , with a notion of distance, d ,
 890 between points). A mapping $\mathcal{T}: X \rightarrow X$ is a contraction mapping if there exists a constant $c: 0 \geq c > 1$
 891 such that for all $x \in X$:

$$892 \quad d(\mathcal{T}[x_i], \mathcal{T}[x_j]) \leq cd(x_i, x_j)$$

893 That is, a contraction mapping maps points closer together.

894 Elaborating now on the operator $\mathcal{T}_{\alpha\tau}[V](s)$ and using $|\cdot|$ as our distance metric:

$$895 \quad |\mathcal{T}_{\alpha K}[V_1](i) - \mathcal{T}_{\alpha K}[V_2](i)|$$

$$896 \quad = \left| V_1(i) + \alpha \cdot \sum_{i,j \in \mathcal{S}} p_{ij} \cdot \mathcal{X}^K \cdot (r_{i,j} + \gamma V_1(j) - V_1(i)) - V_2(i) + \alpha \right.$$

$$897 \quad \left. \cdot \sum_{i,j \in \mathcal{S}} p_{ij} \cdot \mathcal{X}^K \cdot (r_{i,j} + \gamma V_2(j) - V_2(i)) \right|$$

$$\begin{aligned}
 898 \quad & |\mathcal{T}_{\alpha K}[V_1](i) - \mathcal{T}_{\alpha K}[V_2](i)| \\
 899 \quad & = \left| V_1(i) - V_2(i) \right. \\
 900 \quad & \left. + \alpha \sum_{i,j \in S} p_{ij} \cdot \left(\mathcal{X}^\tau \cdot (r_{i,j} + \gamma V_1(j) - V_1(i)) - \mathcal{X}^\tau \cdot (r_{i,j} + \gamma V_2(j) - V_2(i)) \right) \right|
 \end{aligned}$$

901 Using the relation defined above $\mathcal{E}_{a,b,K} \cdot (a - b) = \mathcal{X}^K(a) - \mathcal{X}^K(b)$

$$\begin{aligned}
 902 \quad & |\mathcal{T}_{\alpha K}[V_1](i) - \mathcal{T}_{\alpha K}[V_2](i)| \\
 903 \quad & = \left| V_1(i) - V_2(i) + \alpha \sum_{i,j \in S} p_{ij} \cdot \mathcal{E}_{V_1, V_2, K} \cdot (\gamma(V_1(j) - V_2(i)) - (V_1(j) - V_2(i))) \right|
 \end{aligned}$$

$$\begin{aligned}
 904 \quad & |\mathcal{T}_{\alpha K}[V_1](i) - \mathcal{T}_{\alpha K}[V_2](i)| \\
 905 \quad & = \left| \left(1 - \alpha \sum_{i,j \in S} p_{ij} \cdot \mathcal{E}_{V_1, V_2, K} \right) \cdot (V_1(j) - V_2(i)) \right. \\
 906 \quad & \left. + \sum_{i,j \in S} p_{ij} \cdot \mathcal{E}_{V_1, V_2, K} \cdot (\gamma(V_1(j) - V_2(i))) \right|
 \end{aligned}$$

907

908 Given that $\mathcal{E}_{a,b,K} \in [1 - |K|, 1 + |K|]$ and assuming $\alpha \in (0, (1 + |K|)^{-1})$:

$$909 \quad 1 - \alpha \cdot \sum_{i,j \in S} p_{ij} \cdot \mathcal{E}_{V_1, V_2, K} > 0$$

910 Taking this term outside the $|\cdot|$ and rearranging:

$$911 \quad |\mathcal{T}_{\alpha K}[V_1](i) - \mathcal{T}_{\alpha K}[V_2](i)| = \left(1 - \alpha \cdot (1 - \gamma) \cdot \sum_{i,j \in S} p_{ij} \cdot \mathcal{E}_{V_1, V_2, K} \right) |(V_1(j) - V_2(i))|$$

$$912 \quad |\mathcal{T}_{\alpha K}[V_1](i) - \mathcal{T}_{\alpha K}[V_2](i)| = c \cdot |(V_1(j) - V_2(i))|$$

913 Where the term:

$$914 \quad c = \left(1 - \alpha \cdot (1 - \gamma) \cdot \sum_{i,j \in S} p_{ij} \cdot \mathcal{E}_{V_1, V_2, K} \right)$$

915 To get the upper boundary of c we use the minimum value for the sum, where $\mathcal{E}_{V_1 V_2 K} = 1 - |K| \quad \forall i, j \in$
 916 S . And use the assumption that $\alpha \in (0, (1 + |K|)^{-1})$:

$$917 \quad c \leq (1 - \alpha \cdot (1 - \gamma) \cdot (1 - |K|))$$

$$918 \quad \lim_{\alpha \rightarrow 0} c = (1 - \alpha \cdot (1 - \gamma) \cdot (1 - |K|)) = 1$$

919 To get the lower boundary of c we use the maximum value for the sum, where $\mathcal{E}_{V_1 V_2 K} = 1 + |K| \quad \forall i, j \in$
 920 S . And use the assumption that $\alpha \in (0, (1 + |K|)^{-1})$:

$$921 \quad c \geq (1 - \alpha \cdot (1 - \gamma) \cdot (1 + |K|))$$

$$922 \quad \lim_{\alpha \rightarrow (1+|K|)^{-1}} c = (1 - (1 + |K|)^{-1} \cdot (1 - \gamma) \cdot (1 + |K|)) = \gamma$$

923 Therefore: $\gamma < c < 1$

$$924 \quad |\mathcal{T}_{\alpha K}[V_1](i) - \mathcal{T}_{\alpha K}[V_2](i)| \leq c \cdot |(V_1(j) - V_2(i))|$$

925 And the operator $\mathcal{T}_{\alpha \tau}[V](s)$ is a contraction mapping, under the condition $\alpha \in (0, (1 + |K|)^{-1})$

926 **1.3.2. Convergence of TD learning with D1 and D2 populations**

927 We define the D1-D2 TD-learning rule as:

$$928 \quad \hat{V}_n(s) \leftarrow \hat{V}_{n-1}(s) + \alpha \cdot \mathcal{X}^\tau(\delta_{s_{n-1}, s_n}) - \beta \hat{V}_{n-1}(s)$$

929 Note this update rule is analogous to the risk sensitive TD learning rule except for the last term that
 930 emerges from the decay factor in the P, N populations of our model.

931 Performing the same re-arrangement as above we reach:

$$932 \quad \hat{V}_n(s) \leftarrow (1 - \sigma/\alpha) \hat{V}_{n-1}(s) + \sigma/\alpha (\alpha \cdot \mathcal{X}^\tau(\delta_{n-1}) + \hat{V}_{n-1}(s) - \alpha \cdot \beta \hat{V}_{n-1}(s))$$

933 We define a new operator $\mathcal{T}'_{\alpha K}[V](s)$:

$$934 \quad \mathcal{T}'_{\alpha K}[V](s) := (1 - \alpha \cdot \beta) \cdot V(s) + \alpha \cdot \sum_{i,j \in S} p_{ij} \cdot \mathcal{X}^K \cdot \delta_{ij}$$

935 And the noise then is defined as:

$$936 \quad \omega_{n-1}(s) = \alpha \cdot \mathcal{X}^K(\delta_{s_{n-1}, s_n}) + \hat{V}_{n-1}(s) - \alpha \cdot \beta \cdot \hat{V}_{n-1}(s) - \mathcal{T}'_{\alpha K}[\hat{V}_{n-1}](s)$$

937 The update becomes:

$$938 \quad \hat{V}_n(s) \leftarrow (1 - \sigma/\alpha)\hat{V}_{n-1}(s) + \sigma/\alpha \left(\mathcal{T}'_{\alpha K}[\hat{V}_{n-1}](s) + \omega_{n-1}(s) \right)$$

939 The noise term reduces to the same expression as the one of TD learning, and so it fullfils the
 940 requirements for the theorem of stochastic iterative algorithms. We will now test whether the operator
 941 $\mathcal{T}'_{\alpha K}[V](s)$ also represents a contraction map.

$$942 \quad | \mathcal{T}_{\alpha K}[V_1](i) - \mathcal{T}_{\alpha K}[V_2](i) |$$

$$943 \quad = \left| V_1(i) + \alpha \cdot \sum_{i,j \in \mathcal{S}} p_{ij} \cdot \mathcal{X}^K \cdot (r_{i,j} + \gamma V_1(j) - V_1(i)) - \alpha \cdot \beta \cdot V_1(i) + V_2(i) + \alpha \right.$$

$$944 \quad \left. \cdot \sum_{i,j \in \mathcal{S}} p_{ij} \cdot \mathcal{X}^K \cdot (r_{i,j} + \gamma V_2(j) - V_2(i)) - \alpha \cdot \beta \cdot V_2(i) \right|$$

$$945 \quad | \mathcal{T}_{\alpha K}[V_1](i) - \mathcal{T}_{\alpha K}[V_2](i) |$$

$$946 \quad = \left| V_1(i) - V_2(i) - \alpha \cdot \beta (V_1(i) - V_2(i)) \right.$$

$$947 \quad \left. + \alpha \sum_{i,j \in \mathcal{S}} p_{ij} \cdot \left(\mathcal{X}^\tau \cdot (r_{i,j} + \gamma V_1(j) - V_1(i)) - \mathcal{X}^\tau \cdot (r_{i,j} + \gamma V_2(j) - V_2(i)) \right) \right|$$

948 Using the relation defined above $\mathcal{E}_{a,b,K} \cdot (a - b) = \mathcal{X}^K(a) - \mathcal{X}^K(b)$

$$949 \quad | \mathcal{T}_{\alpha K}[V_1](i) - \mathcal{T}_{\alpha K}[V_2](i) |$$

$$950 \quad = \left| (1 - \alpha\beta) \cdot (V_1(i) - V_2(i)) \right.$$

$$951 \quad \left. + \alpha \sum_{i,j \in \mathcal{S}} p_{ij} \cdot \mathcal{E}_{V_1, V_2, K} \cdot \left(\gamma (V_1(j) - V_2(i)) - (V_1(j) - V_2(i)) \right) \right|$$

$$\begin{aligned}
 & |\mathcal{T}_{\alpha K}[V_1](i) - \mathcal{T}_{\alpha K}[V_2](i)| \\
 &= \left| \left(1 - \alpha \cdot \beta - \alpha \sum_{i,j \in S} p_{ij} \cdot \mathcal{E}_{V_1, V_2, K} + \gamma \alpha \sum_{i,j \in S} p_{ij} \cdot \mathcal{E}_{V_1, V_2, K} \right) \cdot (V_1(j) - V_2(i)) \right|
 \end{aligned}$$

$$|\mathcal{T}_{\alpha K}[V_1](i) - \mathcal{T}_{\alpha K}[V_2](i)| = c \cdot |(V_1(j) - V_2(i))|$$

Where the term:

$$c = \left(1 - \alpha \cdot \beta - \alpha \cdot (1 - \gamma) \cdot \sum_{i,j \in S} p_{ij} \cdot \mathcal{E}_{V_1, V_2, K} \right)$$

To get the upper boundary of c we use the minimum value for the sum, where $\mathcal{E}_{V_1, V_2, K} = 1 - |K| \quad \forall i, j \in S$, and use the assumption that $\alpha \in (0, (1 + |K|)^{-1})$:

$$c \leq 1 - \alpha \cdot \beta - \alpha \cdot (1 - \gamma) \cdot (1 - |K|)$$

$$\lim_{\alpha \rightarrow 0} c = 1 - \alpha \cdot \beta - \alpha \cdot (1 - \gamma) \cdot (1 - |K|) = 1$$

To get the lower boundary of c we use the maximum value for the sum, where $\mathcal{E}_{V_1, V_2, K} = 1 + |K| \quad \forall i, j \in S$, and use the assumption that $\alpha \in (0, (1 + |K|)^{-1})$:

$$c \geq (1 - \alpha \cdot \beta - \alpha \cdot (1 - \gamma) \cdot (1 + |K|))$$

$$\lim_{\alpha \rightarrow (1+|K|)^{-1}} c = (1 - (1 + |K|)^{-1} \cdot \beta - (1 + |K|)^{-1} \cdot (1 - \gamma) \cdot (1 + |K|)) = \gamma - (1 + |K|)^{-1} \cdot \beta$$

Given that we want $c \geq 0$ we can find the parameter ranges to achieve this:

$$c = \gamma - (1 + |K|)^{-1} \cdot \beta \geq 0$$

Given that: $K \in (-1, 1)$, we use the minimum value of $|K| = 0$ to find the limit of c :

$$\lim_{|K| \rightarrow 0} c = \gamma - \beta$$

So the condition $\gamma \geq \beta$ needs to be present to keep: $0 \leq c < 1$.

Under these conditions, the operator $\mathcal{T}'_{\alpha K}[V](s)$ also represents a contraction map.

$$|\mathcal{T}'_{\alpha K}[V_1](i) - \mathcal{T}'_{\alpha K}[V_2](i)| \leq c \cdot |(V_1(j) - V_2(i))|$$

Stochastic fixed point for the value estimate:

973 Having shown convergence of the algorithm, we will now derive the convergent points for our algorithm
 974 using stochastic fixed points. For clarity, we estimate the stochastic fixed point dropping the dependency
 975 on time.

976 If learning between P and N is symmetric $\alpha^+ = \alpha^- = \alpha$. We derive the convergent estimate of $V(s_t)$
 977 with a fixed-point approach. First, we subtract the P and N update equations, to arrive to the update in
 978 the $\hat{V}(s_t)$ between the (n) and the $(n + 1)$ update:

$$979 \quad \hat{V}^{(n+1)}(s_i) \leftarrow \hat{V}^{(n)}(s_i) + \alpha \cdot \delta^{(n)} - \beta \cdot \hat{V}^{(n)}(s_i)$$

980 Where the superscripts indicate the iteration number. We can now derive the stochastic fixed-point
 981 for $\hat{V}(s_t)$:

$$982 \quad \mathbf{E}[\hat{V}^{(n+1)}(s_i) - \hat{V}^{(n)}(s_i)] = 0$$

$$983 \quad \mathbf{E}[\alpha \cdot \delta^{(n)} - \beta \cdot \hat{V}^{(n)}(s_i)] = 0$$

$$984 \quad \mathbf{E}[\alpha \cdot (r^{(n)} - \hat{V}(s_i)) - \beta \cdot \hat{V}^{(n)}(s_i)] = 0$$

$$985 \quad \alpha \cdot \mathbf{E}[r] - (\alpha + \beta) \cdot \mathbf{E}[\hat{V}(s_i)] = 0$$

$$986 \quad \mathbf{E}[\hat{V}(s_i)] = \frac{\alpha}{\alpha + \beta} \mathbf{E}[r]$$

$$987 \quad \mathbf{V}^* = \frac{\alpha}{\alpha + \beta} \mathbf{E}[r]$$

988 Where $\mathbf{V}^* = \mathbf{E}[\hat{V}(s_i)]$ is the value around which $\hat{V}(s_i)$ is expected to fluctuate after convergence.

989 Throughout this study, we have manipulated the learning rates between P and N to be asymmetric $\alpha^+ \neq$
 990 α^- or, equivalently, $\tau \neq 1 - \tau$. We can find the stochastic fixed point for this more general case:

$$991 \quad \mathbf{E}[\hat{V}^{(n+1)}(s_i) - \hat{V}^{(n)}(s_i)] = 0$$

$$992 \quad \mathbf{E}[\tau \cdot |\delta^{(n)}| \cdot \mathbf{1}_{\delta > 0} - (1 - \tau) \cdot |\delta^{(n)}| \cdot \mathbf{1}_{\delta < 0} - \beta \cdot \hat{V}^{(n)}(s_i)] = 0$$

993 To take the expectation we use the definition: $\mathbf{E}[X] = \sum_i p(x_i) \cdot x_i$. For a Bernoulli distribution, $p(x_i)$
 994 takes two values:

- 995 • $p(x_i) = p$ if reward is delivered and, thus $r = 1$, $\delta_t > 0$,
- 996 • $p(x_i) = (1 - p)$ if reward is not delivered and, thus $r = 0$, $\delta_t < 0$,

997 Therefore, we can resolve the expectation and expand the RPEs:

$$998 \quad \mathbf{E}[\tau \cdot |\delta^{(n)}| \cdot \mathbf{1}_{\delta > 0} - (1 - \tau) \cdot |\delta^{(n)}| \cdot \mathbf{1}_{\delta < 0} - \beta \cdot \hat{V}^{(n)}(s_i)] = 0$$

$$999 \quad \tau \cdot \mathbf{E}[|r - \hat{V}(s_i)| \cdot \mathbf{1}_{\delta > 0}] - (1 - \tau) \cdot \mathbf{E}[|-\hat{V}(s_i)| \cdot \mathbf{1}_{\delta < 0}] - \beta \cdot \mathbf{E}[\hat{V}(s_i)] = 0$$

1000 Taking the absolute values:

$$1001 \quad |r - \hat{V}(s_i)| = r - \hat{V}(s_i) \dots \text{ if } (r - \hat{V}(s_i)) > 0$$

$$1002 \quad |-\hat{V}(s_i)| = \hat{V}(s_i) \dots \text{ if } (-\hat{V}(s_i)) < 0$$

$$1003 \quad \tau \cdot \mathbf{E}[(r - \hat{V}(s_i)) \cdot \mathbf{1}_{\delta > 0}] - (1 - \tau) \cdot \mathbf{E}[\hat{V}(s_i) \cdot \mathbf{1}_{\delta < 0}] - \beta \cdot \mathbf{E}[\hat{V}(s_i)] = 0$$

1004 Replacing stochastic fixed point: $\mathbf{E}[\hat{V}(s_t)] = V^*$ and taking the expectations:

$$1005 \quad \beta \cdot V^* = \tau \cdot p \cdot (r - V^*) - (1 - \tau) \cdot (1 - p) \cdot V^*$$

1006 Rearranging and isolating V^* , we obtain:

$$1007 \quad V^* = \frac{\frac{\tau}{1 - \tau} \cdot \frac{p}{1 - p} \cdot r}{\frac{\tau}{1 - \tau} \cdot \frac{p}{1 - p} + 1 + \frac{\beta}{(1 - \tau) \cdot (1 - p)}}$$

1008 *Stochastic fixed point for P and N populations:*

1009 We have mentioned that the decay term (β) in the update equations serves to stabilize the estimates of the
1010 P and N populations (i.e., avoid infinite increases). We can observe the influence of β by computing the
1011 stochastic fixed points for these variables.

1012 For the P population:

$$1013 \quad \mathbf{E}[P^{(n+1)}(s_t) - P^{(n)}(s_t)] = 0$$

$$1014 \quad \mathbf{E}[\tau \cdot |r - V^{(n)}(s_t)| \cdot \mathbf{1}_{\delta > 0} - \beta \cdot P^{(n)}(s_t)] = 0$$

$$1015 \quad p \cdot \tau \cdot (r - V^*) - \beta \cdot P^* = 0$$

$$1016 \quad P^* = \frac{p \cdot \tau}{\beta} \cdot (r - V^*)$$

1017 Similarly, for the N population:

1018
$$\mathbf{E}[N^{(n+1)}(s_t) - N^{(n)}(s_t)] = 0$$

1019
$$\mathbf{E}[(1 - \tau) \cdot |-V^{(n)}| \cdot \mathbf{1}_{\delta < 0} - \beta \cdot N^{(n)}(s_t)] = 0$$

1020
$$(1 - p) \cdot (1 - \tau) \cdot V^* - \beta \cdot N^* = 0$$

1021
$$N^* = \frac{(1 - p) \cdot (1 - \tau)}{\beta} \cdot V^*$$

1022 As it can be seen in the stochastic fixed points P^*, N^* , the term $\frac{1}{\beta}$ is a proportionality constant. Therefore:

1023
$$\lim_{\beta \rightarrow 0} P^* = \lim_{\beta \rightarrow 0} \left(\frac{p \cdot \tau}{\beta} \cdot (r - V^*) \right) = \text{undefined}$$

1024
$$\lim_{\beta \rightarrow 0} N^* = \lim_{\beta \rightarrow 0} \left(\frac{(1 - p) \cdot (1 - \tau)}{\beta} \cdot V^* \right) = \text{undefined}$$

1025 So, $\beta \neq 0$ needs to be met for the stochastic fixed points P^*, N^* to exist. In Extended Data Fig. 1 we show
1026 empirically that the convergence rate is slower as β gets closer to 0, but it is always achieved.

1027 **1.3.3. Sensitivity of learned variables in D1-D2 model to parameters**

1028 The conditions for the D1-D2 model to reproduce the data from our habenula lesion experiment and some
1029 of the previous studies are that:

- 1030 1. The bias in V^* induced by the asymmetric learning rates doesn't change the monotonicity of the
1031 learned values as a function of the true expected value of the return distribution $\mathbf{E}[R(s)]$. In other
1032 words, regardless of the level of 'optimism' or 'pessimism', V^* monotonically increases with
1033 $\mathbf{E}[R(s)]$.
- 1034 2. Asymmetric learning rates change the concavity of V^* as a function of $\mathbf{E}[R(s)]$: 'Optimistic' or
1035 'pessimistic' value functions are concave or convex with respect to $\mathbf{E}[R(s)]$, respectively.

1036

1037 We will now analyze whether these conditions are met, considering the range of parameters of relevance:

1038 $0 < \tau < 1, r \neq 0$ and $0 < \beta < 1$

1039 For the condition 1 to be met, the first derivative of V^* with respect to $\mathbf{E}[R(s)]$ should always be positive.

1040 In the case of Bernoulli return distributions, the derivative of V^* with respect to $p(\text{reward})$ is

$$1041 \quad \frac{\partial V^*}{\partial p} = \frac{\partial}{\partial p} \left(\frac{\frac{\tau}{1-\tau} \cdot \frac{p}{1-p} \cdot r}{\frac{\tau}{1-\tau} \cdot \frac{p}{1-p} + 1 + \frac{\beta}{(1-\tau) \cdot (1-p)}} \right) = \frac{\tau \cdot r \cdot (\tau - \beta - 1)}{(\tau \cdot (2p - 1) + \beta - \tau + 1)^2}$$

1042 We can look at the fixed points of this expression, as they correspond to the value of p at which the
 1043 derivative changes the sign. This expression has fixed points at: $\tau = 0$, $r = 0$, and $\beta = \tau - 1$. Given our
 1044 parameters' ranges: $0 < \tau < 1$, $r \neq 0$ and $0 < \beta < 1$, none of those fixed points are present within those
 1045 ranges. In addition, it can be seen that this $\left(\frac{\partial V^*}{\partial p}\right)$ is positive for the parameter values within those ranges.

1046 Therefore, knowing that the derivative won't reach any fixed point, V^* is always a growing monotonic
 1047 function with respect to p .

1048 For the condition 2 to be met, we can analyze the second derivative of V^* with respect to $\mathbf{E}[R(s)]$ as it
 1049 indicates the convexity of a function. The conditions to be met are:

$$1050 \quad - \quad V^* \text{ is convex if it is 'pessimistic': if } \tau < 0.5 \rightarrow \frac{\partial^2 V}{\partial p^2} > 0$$

$$1051 \quad - \quad V^* \text{ is concave if it is 'optimistic': if } \tau > 0.5 \rightarrow \frac{\partial^2 V}{\partial p^2} < 0$$

1052 In the case of Bernoulli return distributions, we take the second derivative of V^* with respect to
 1053 p (reward):

$$1054 \quad \frac{\partial^2 V^*}{\partial p^2} = \frac{\partial^2}{\partial p^2} \left(\frac{\frac{\tau}{1-\tau} \cdot \frac{p}{1-p} \cdot r}{\frac{\tau}{1-\tau} \cdot \frac{p}{1-p} + 1 + \frac{\beta}{(1-\tau) \cdot (1-p)}} \right) = \frac{2\tau \cdot (2\tau - 1) \cdot r \cdot (\tau - \beta - 1)}{(\tau \cdot (2p - 1) + \beta - p + 1)^3}$$

1055 We can again look at the fixed points of this expression. These happen at: $\tau = 0$, $r = 0$, and $\beta = \tau - 1$
 1056 and $\tau = 0.5$. Among them, the only fixed point within our parameters range is the latter. In addition, by
 1057 replacing τ in the expression above, it is easily shown that it is positive if $\tau < 0.5$ and negative if $\tau > 0.5$.
 1058 Thus, given that the ranges for the parameters are such that the second derivative won't reach any other
 1059 fixed point, condition 2 will always be met.

1060 **1.4. Distributional TD learning with D1 and D2 populations**

1061 The signatures of distributional RL were preserved in dopamine neurons firing rates after habenula
 1062 lesions (Extended Data Fig. 3-4). Therefore, we considered a third alternative to model 1 and 2, that
 1063 assigns different functions to each of the mechanisms for asymmetric learning rates.

1064 In this model (Extended Data Fig. 13) the single cell asymmetric scaling factors (α_i^+ , α_i^-) give rise to a
 1065 distributional expectile code for value and are implemented at the level of the scaling of RPE-evoked
 1066 responses of dopamine neurons:

1067
$$\delta_{i,t} := r_t + \gamma \cdot \bar{z}(s_{t+1}) - V_i(s_t)$$

1068
$$\hat{\delta}_{i,t} = \alpha_i^+ \cdot \delta_{i,t} \dots \text{if } \delta_{i,t} > 0$$

1069
$$\hat{\delta}_{i,t} = \alpha_i^- \cdot \delta_{i,t} \dots \text{if } \delta_{i,t} < 0$$

1070

1071 The modulation of receptor sensitivities, carried out downstream at the SPN level, gives rise to the global
 1072 rescaling of the value updates (η^+ , η^-) (Extended Data Fig. 13A):

1073
$$P_i(s_t) \leftarrow P_i(s_t) + \eta^+ \cdot |\hat{\delta}_i(t)| - \beta \cdot P_i(s_t) \dots \text{if } \delta_i(t) > 0$$

1074
$$N_i(s_t) \leftarrow N_i(s_t) + \eta^- \cdot |\hat{\delta}_i(t)| - \beta \cdot N_i(s_t) \dots \text{if } \delta_i(t) > 0$$

1075
$$\hat{V}_i(s_t) = P_i(s_t) - N_i(s_t)$$

1076 These set of update equations are equivalent to a modified version of the update equation of distributional
 1077 RL:

1078
$$[\Delta V_i(s_t)] = \frac{1}{N} \sum_j^N \eta^+ \cdot \alpha_i^+ \cdot \delta_{i,j} \cdot I_{\delta_{i,j} > 0} + \eta^- \cdot \alpha_i^- \cdot \delta_{i,j} \cdot I_{\delta_{i,j} > 0}$$

1079
$$V_i(s_t) \leftarrow V_i(s_t) + \mathbf{E}[\Delta V_i(s_t)]$$

1080 Thus, this model can give rise to biases in value learning (Extended Data Fig. 13), while keeping intact
 1081 information about the value distribution. By employing the results from the biophysical model (Fig. 6),
 1082 we found that this distributional TD model can parsimoniously explain all aspects of the data in the
 1083 habenula lesion study (Extended Data Fig. 12B), including the features of a distributional code and the
 1084 optimistic biases observed in behavior and dopamine cue-evoked responses (Extended Data Fig. 12).

1085 **1.5. Dependency of model on assumption: Log vs. linear scaling of receptor occupancy curves**

1086 Through this work, we have used the dose-occupancy curves of D1 and D2 receptors to derive the
1087 receptor sensitivities that result in the asymmetric scaling factors in Model 1. It is important to note that
1088 the slopes of the receptor occupancy curve (= receptor sensitivity) were obtained from the receptor
1089 occupancy curves plotted as a function of log of dopamine concentrations.

$$1090 \quad \alpha^+ = \frac{\Delta O_{ccD1}}{\Delta \log(C_{DA^+})}$$

$$1091 \quad \alpha^- = \frac{\Delta O_{ccD2}}{\Delta \log(C_{DA^-})}$$

1092

1093 To show that this assumption is not essential, we now derive the receptors sensitivities assuming linear
1094 changes in dopamine levels due to RPE-evoked responses.

$$1095 \quad \alpha^+ = \frac{\Delta O_{ccD1}}{\Delta C_{DA^+}}$$

$$1096 \quad \alpha^- = \frac{\Delta O_{ccD2}}{\Delta}$$

1097 As shown in Extended Data Fig. 9, the choice of a linear versus log scale affects the absolute magnitude
1098 of the derived receptor sensitivities, but the normalized metric $\tau = \frac{\alpha^+}{\alpha^- + \alpha^+}$ holds the same relationship to
1099 baseline dopamine levels with a small shift in the curve (Extended Data Fig.9, right panel). The
1100 normalized metric is the factor determining the update asymmetries and, thus, the stochastic fixed points
1101 at which the variables converge.

1102 **1.6. Normative motivation for two-factor learning rule**

1103 We have used in the previous models a so-called *two factor learning rule*, where the value updates
1104 depend only on the presynaptic activity (i.e., state input) and TD RPEs. Here, we motivate this choice
1105 from a normative approach based on previous work¹¹.

1106 Consider a linear approximation for value, where the value function (\hat{V}) is the output of a single linear
 1107 neuron. Here, \hat{V} is a linear function of the input feature-vector representing the state $\mathbf{x}(s) =$
 1108 $(x_1(s), \dots, x_n(s))$, parametrized with a weight vector $\mathbf{w} = (w_1, \dots, w_n) . . :$

$$1109 \quad \hat{V}(s, \mathbf{w}) = \mathbf{w}^T \mathbf{x}(s)$$

1110 To put it into neural terms, we can think of $x_i(s)$ as the presynaptic activity onto the value neuron \hat{V} ,
 1111 with a synaptic efficacy w_i .

1112 As before, the agent computes the TD error based on this linear approximation for value and the sampled
 1113 reward:

$$1114 \quad \delta_t = r_t + \gamma \cdot \hat{V}(s_{t+1}, \mathbf{w}) - \hat{V}(s_t, \mathbf{w})$$

1115 In the problem of *value prediction*, the agent aims to achieve the highest accuracy of prediction. One way
 1116 to achieve this is to perform *stochastic gradient descent* (SGD) with respect to the parameters (\mathbf{w}) of the
 1117 value function to minimize the objective function such as the squared error (δ_t^2). We can define this
 1118 optimization problem as: $\operatorname{argmin}_{\mathbf{w}} \left(\frac{1}{2} \delta^2 \right)$ where we have deliberately chosen the constant $\frac{1}{2}$ for clarity,
 1119 but it doesn't change the end results.

1120 To perform SGD in this minimization problem, the parameters (\mathbf{w}) should be updated in the opposite
 1121 direction of the gradient of the loss with respect to the parameters (i.e., opposite to $\nabla_{\mathbf{w}} \left(\frac{1}{2} \delta^2 \right)$):

$$1122 \quad \mathbf{w} \leftarrow \mathbf{w} - \alpha \cdot \nabla_{\mathbf{w}} \left(\frac{1}{2} \delta^2 \right)$$

1123 Where α is the learning rate. To compute the gradient, we use the chain rule:

$$1124 \quad \nabla_{\mathbf{w}} \left(\frac{1}{2} \delta^2 \right) = \frac{\partial (1/2 \delta^2)}{\partial \hat{V}} \cdot \frac{\partial \hat{V}}{\partial \mathbf{w}} = \frac{\partial (1/2 (r - \hat{V})^2)}{\partial \hat{V}} \cdot \frac{\partial \mathbf{w}^T \mathbf{x}(s_t)}{\partial \mathbf{w}} = \frac{-2(r - V)}{2} \cdot \mathbf{x}(s_t) = -\delta \cdot \mathbf{x}(s_t)$$

1125 Therefore, the update for the parameters of the value function is:

$$1126 \quad \mathbf{w} \leftarrow \mathbf{w} + \alpha \cdot \delta \cdot \mathbf{x}(s_t)$$

1127 The term $\delta \cdot \mathbf{x}(s_t)$ in the equation above is what we call a *two-factor learning rule*, dependent only on the
 1128 presynaptic activity and not contingent on the post-synaptic activity.

1129 The development of the TD learning model with D1 and D2 populations (section 1.3) has respected this
 1130 learning rule, complying with what is required for SGD in the value prediction problem. Note that we

1131 have implicitly developed our models with a complete serial compound representation (CSC) of the
1132 states¹¹, where $x(s_t) = 1$ in a single element $x_i(s_t)$ representing the current state and 0 otherwise. It can
1133 be shown that with this representation, the update equation above is equivalent to:

$$1134 \quad \hat{V} \leftarrow \hat{V} + \alpha \cdot \delta$$

1135

1136 **2. Computational model of dopamine release and receptor occupancy**

1137 To predict changes in dopamine concentrations and receptor occupancies (Fig. 6), we employed a
1138 biophysical model developed elsewhere⁵⁹. It presents two interacting dynamical systems. The first
1139 system models the change in receptor occupancies while the second the change in dopamine levels per
1140 unit time.

1141 In the first system, the occupancy of receptors is modelled as a binding reaction between dopamine (*DA*)
1142 and D1 or D2 receptors (*R*), using the constants for forward and backward reactions (k_{on} , k_{off}).



1144 This formulation results in the following equation for the change in receptor occupancy $Occ(t)$ per unit
1145 time:

$$1146 \quad \frac{dOcc(t)}{dt} = (1 - Occ(t)) \times k_{on} \times C_{DA}(t) - Occ(t) \times k_{off}$$

1147 The values used for the association and dissociation constants for each receptor type (k_{on} and k_{off} ,
1148 respectively) are detailed in Table 1.

1149 In the second system, the change in dopamine concentration ($C_{DA}(t)$) is a function of both dopamine
1150 release and uptake.

$$1151 \quad \frac{dC_{DA}(t)}{dt} = DA_{release}(t) - DA_{uptake}(t)$$

1152 Dopamine release is a product of firing rate ($v(t)$) and release capacity ($\gamma(t)$)

$$1154 \quad DA_{release}(t) = \gamma(t) \cdot v(t)$$

1153 Where:

- 1155 1. $v(t)$ is the firing rate of dopamine neurons, provided by the neural data.
- 1156 2. $\gamma(t) = \gamma_{pr_n} \cdot P_r \cdot G_{D2}(t)$ is defined as the increase in $C_{DA}(t)$ by a single synchronized action
 1157 potential:
- 1158 a. $P_r = 1$ release probability in the absence of presynaptic D2-autorreceptors,
- 1159 b. $\gamma_{pr_n} = 2$ release capacity in the absence of presynaptic D2-autorreceptors. This value was set
 1160 to be deliberately high and anticipates a ~50% reduction by terminal feedback.
- 1161 c. $G_{D2}(t)$ is a multiplicative gain that represents the modulation of dopamine release by D2-
 1162 autorreceptors. This is a decaying function of the occupancy of D2-autorreceptors
 1163 ($Occ_{D2_a}(t)$) which is modelled by the same binding reaction explained above. The gain is
 1164 parametrized by the autoreceptor efficacy, $\alpha = 3$. The smaller the α the less the decay in
 1165 release with receptor occupancy.

$$G_{D2}(t) = \frac{1}{1 + \alpha \cdot Occ_{D2_a}(t)}$$

1167 Dopamine uptake is a function of the uptake of dopamine by the dopamine transporter (DAT) and other
 1168 non-DAT sources

$$1169 \quad DA_{uptake}(t) = dt \cdot \left(\frac{V_{max}^{pr_n} \cdot C_{DA}(t)}{K_m + C_{DA}(t)} - K_{nonDAT} \right)$$

1170 Where:

- 1171 - $V_{max}^{pr_n} = 1500 \frac{nM}{sec}$ the maximal uptake capacity assuming approximately 100 terminals in the
 1172 near surroundings.
- 1173 - $K_m = 160 nM$, is the Michaelis-Menten parameter for uptake mediated by DAT
- 1174 - $K_{nonDAT} = 0.04 nM$ is a constant for the dopamine removal not mediated by DAT. For
 1175 example, monoamine oxidase (MAO) and noepinephrine transporter (NET) mediated uptake.

1176 The variables of the model reported in Fig. 6 correspond to: $Occ_{D1}(t)$, $Occ_{D2}(t)$, $C_{DA}(t)$. We used as
 1177 input to the model the firing rates derived from the electrophysiological recording of optogenetically
 1178 identified dopamine neurons conducted in Tian and Uchida (2015)⁵⁵. This modeling, while considering
 1179 major processes, does not take into account all of the complexity of the biological environment in the

1180 brain, yet we used this model to obtain an approximate estimate of the order of changes in dopamine
1181 concentrations and receptor occupancies.

1182 3. Simulation details of habenula lesion data

1183 3.1. Biophysical model simulations

1184 We used the computational model described previously (methods section 2)⁵⁹ to calculate the
1185 extracellular dopamine levels and estimate the occupancy of postsynaptic receptors from the habenula
1186 lesion dataset. The model was driven by the average spike rate of dopamine neurons recorded from
1187 control or lesioned animals. For each recorded dopamine neuron, the simulations were carried on a trial
1188 by trial basis that consisted of a time window [-15, 20] sec with respect to cue onset. A relatively large
1189 window was used to allow for the relevant variables to stabilize in its baseline, as the simulations were
1190 initialized at zero.

1191 For each trial, spikes were first binned with 10-ms windows and then smoothed by a Gaussian kernel
1192 ($\sigma = 0.3 \times (ISI_{mean})$). All trials were then averaged across trials, to determine the mean single-cell
1193 response for dopamine release and D1 and D2 receptor activation. Final average dopamine concentrations
1194 and receptor occupancies were obtained from the average of all mean single-cell responses.

1195 *Computation of receptors sensitivities from the model results*

1196 We computed the receptor sensitivity from the occupancies Occ_{D1} , Occ_{D2} and their theoretical dose-
1197 occupancy curves. Starting from the occupancy at baseline, we derived the change in occupancy as a
1198 function of the transients in dopamine concentration C_{DA} elicited by RPE-evoked dopamine responses, at
1199 the level of the population average.

1200 The ratio between these quantities corresponds to the receptors' sensitivities. These are transferred as α^+
1201 and α^- to our reinforcement learning model (model 1):

$$1202 \quad \alpha^+ = \frac{\Delta Occ_{D1}}{\Delta C_{DA}} \dots \text{ if } \Delta C_{DA} > 0$$

$$1203 \quad \alpha^- = \frac{\Delta Occ_{D2}}{\Delta C_{DA}} \dots \text{ if } \Delta C_{DA} < 0$$

1204 Where ΔC_{DA} , ΔOcc_{D1} , ΔOcc_{D2} are the changes computed with respect to baseline, as: $\Delta x =$
1205 $\bar{x}_{outcome} - \bar{x}_b$, for each variable $x = \{C_{DA}, Occ_{D1}, Occ_{D2}\}$. Where \bar{x} denotes the population average

1206 response for each group. The outcome responses were taken as the average from [0,1] sec after outcome
1207 onset, while the baseline was taken as the average from [-1, 0] sec with respect to cue onset.

1208 3.2. Model 1 simulations

1209 The simulations for Model 1 were carried out with a TD learning model with D1 and D2 populations
1210 (methods section 1.3). We ran the simulations using the resultant receptor sensitivities from the
1211 biophysical model as the population-level asymmetric learning rates in Model 1 (i.e., the learning rates
1212 α^+ , α^- for P and N updates). The simulations were run for 3,000 trials on the Pavlovian conditioning task
1213 used in the study⁵⁵. We assumed a uniform distribution of trial types across the session. Each trial
1214 consisted of 4 states (baseline, cue, delay, reward), assuming Markovian dynamics between them. All
1215 variables were initialized at zero. The model had as hyper-parameters a discounting factor of $\gamma = 0.99$
1216 and a decay term $\beta = 0.002$. We report in Fig. 4, Model 1 results assuming a uniform scaling of TD
1217 RPEs across the neuronal population. In Extended Data Fig. 12 we show that this model reproduces key
1218 signatures of the data irrespective of the choice of the decay factor β .

1219 The results are not dependent on a uniform scaling of TD RPEs. Given that distributional RL signatures
1220 were preserved in the data even after habenula lesions, we also considered Model 1 under the
1221 distributional TD learning framework (Extended Data Fig. 13). For this, we used the distribution of single
1222 cell asymmetric scaling factors (α_i^+ , α_i^-) derived from the dopamine neurons firing rates. This model also
1223 reproduced key signatures of the data irrespective of the choice of the decay factor β (Extended Data Fig.
1224 12).

1225 3.3. Model 2 simulations

1226 The simulations for Model 2 were carried out with a TD learning model. As with Model 1, simulations
1227 were run for 3,000 trials on the Pavlovian conditioning task⁵⁵. We assumed a uniform distribution of trial
1228 types across the session. Each trial consisted of 4 states (baseline, cue, delay, reward), assuming
1229 Markovian dynamics between them. All variables were initialized at zero. The model had as parameters a
1230 discounting factor of $\gamma = 0.99$.

1231 We used the distribution of single cell asymmetric scaling factors derived from the firing rates of
1232 dopamine neurons as α_i^+ , α_i^- . In section 1.2 we emphasized that in order to accurately compute the TD
1233 RPE in distributional TD, we require taking samples from the estimated return distribution
1234 $\tilde{z}_i(s_{t+1}) \sim Z(s_{t+1})$. We did this by running an optimization process where we minimize for the expectile

1235 loss between the taken samples $\tilde{z}_i(s_{t+1})$, $V_i(s_{t+1})$ from the model, and τ_i as estimated from the data.

1236 The problem was defined as $\text{argmin}_{s_1 \dots s_m} \mathcal{L}(s, V, \tau)$ where:

1237
$$\mathcal{L}(s, V, \tau) = \frac{1}{M} \sum_{m=1}^M \sum_{i=1}^N |\tau_i - \mathbf{I}_{s_m < V_n}| (\hat{z}_m - V_i)^2$$
, for N neurons and M samples

1238 In the simulations, we took M samples where M equals the number of neurons (N) and performed an

1239 update taking the expectation across all samples as described in the methods section 1.3.

1240 4. Simulations details for replications of previous experimental results

1241 4.1. Cools et al. (2009)

1242 We simulate the results from Cools et al. (2009) (Fig. 7, Extended Data Fig. 6-7) in which they tested the

1243 effects of bromocriptine in altering learning rate asymmetry²⁴. In their study, they performed a reversal

1244 learning task and reported a parameter called ‘relative reversal learning (RRL)’, equivalent to the

1245 difference between the positive and negative learning rates in our model. We computed this as: $\alpha^+ \alpha^+ +$

1246 $+ \alpha^- \alpha^- - \alpha^+ \alpha^- = \tau^+ - 1 - \tau^- = 2\tau - 1$, reported in Fig. 7 E,F, where the parameters α^+ , α^- were

1247 computed from the slopes of the D2l (postsynaptic D2 receptors) and D1 occupancy curves $(2\tau - 1)_{occ}$

1248 or activation curves $(2\tau - 1)_{act}$. The change in relative reversal learning in Fig. 7 H-I was calculated as

1249 taking the difference between drug and the ‘control’ condition as:

1250
$$\Delta(2\tau - 1) = (2\tau - 1)_{drug} - (2\tau - 1)_{control}.$$

1251 We simulated the effect of bromocriptine using the biophysical model for dopamine release and receptor

1252 occupancy (Section 2, Methods). We added an additional ligand for D2 receptors to the update equations

1253 for occupancy:

1254
$$\frac{dOcc_{DA,r_j}(t)}{dt} = \left(1 - Occ_{DA,r_j}(t)\right) \times k_{on}^{DA,r_j} \times C_{DA}(t) - k_{off}^{DA,r_j}$$

1255
$$\frac{dOcc_{Drug,r_j}(t)}{dt} = \left(1 - Occ_{Drug,r_j}(t)\right) \times k_{on}^{Drug,r_j} \times C_{Drug}(t) - k_{off}^{Drug,r_j}$$

1256 Where $r_j: \{D1, D2s, D2l\}$, and $k_{on}^{Drug,D2s} = 0.02083$, $k_{off}^{Drug,D2s} = 0.1$, $k_{on}^{Drug,D2l} = 0.04$, $k_{off}^{Drug,D2l} = 0.1$

1257 are reported in Table 1⁹⁷.

1258 To calculate the effects of efficiency of the drug, we calculated the activation of D2l and D2s receptors in

1259 the following way:

1260
$$Act_{r_j}(t) = E_{DA,r_j} \cdot Occ_{DA,r_j}(t) + E_{Drug,r_j} \cdot Occ_{Drug,r_j}(t)$$

1261 Where $E_{DA,r_j} = 1$ is the efficiency of dopamine on the receptors activation, and $E_{Drug,r_j} < 1$ the
 1262 efficiency of the drug, for $r_j: \{D1, D2s, D2l\}$. The parameter for D1 receptors was kept at $E_{Drug,D1} = 0$
 1263 for all simulations.

1264 To simulate the effects of D2s activation by the drug in D2l occupancy in Fig. 7b,e,h we report the effects
 1265 of $E_{Drug,D2s} = 0$ (solid lines) and $E_{Drug,D2s} = 0.6$ (dashed lines). To simulate the effect of the drug in
 1266 D2s and D2l activation in Fig. 7c,f,i we report the effects of $E_{Drug,D2s} = 0.6, E_{Drug,D2l} = 0.6$.

1267 We show how the qualitative nature of the effects of the drug in relative reversal learning still hold
 1268 regardless of whether the parameter τ is computed from the occupancy curves (Extended Data Fig. 7, Fig.
 1269 7n,e,h) or the activation curves (Extended Data Fig. 8, Fig. 7c,f,i). In addition, in Supplementary Figure
 1270 8-9 we show that the qualitative results still hold regardless of the choice of the efficiency parameters
 1271 $E_{Drug,D2s}$ and $E_{Drug,D2l}$.

1272 **4.2. Timmer et al. (2018)**

1273 In this study²⁵ they reported a ‘loss aversion’ parameter (λ in their results).

1274
$$SUG = (1 - \lambda) \cdot p_{gain} \cdot Gain + \lambda \cdot p_{loss} \cdot Loss$$

1275 Where SUG is the ‘subjective utility’ for a given option, and $p_{gain} = p_{loss}$.

1276 In our formulation, we assume that the task in the study is performed under steady state conditions after
 1277 having learned with a learning rate (τ). With this assumption, the SUG at task performance is equivalent to
 1278 the convergent V estimate after learning. We will show that at these steady state conditions $(1 - \tau)$ is
 1279 equivalent to (λ) .

1280 Starting with the solution for V :

1281
$$SUG = V = \frac{\tau \cdot p_{gain} \cdot r_{gain} + (1 - \tau) \cdot (1 - p_{gain}) \cdot r_{loss}}{\tau \cdot p_{gain} + (1 - \tau) \cdot (1 - p_{gain})}$$

1282 Replacing for $p_{gain} = 0.5$:

1283
$$SUG = \frac{\tau \cdot r_{gain} + (1 - \tau) \cdot r_{loss}}{\tau + (1 - \tau)}$$

1284 Given that: $\tau + (1 - \tau) = 1$

1285
$$\text{SUG} = \tau \cdot r_{\text{gain}} + (1 - \tau) \cdot r_{\text{loss}}$$

1286 Therefore, our model, applied to their task, gives rise to the same SUG computation, with λ equivalent to
1287 $(1 - \tau)$.

1288 To generate Fig. 8F, we performed the following steps:

- 1289 1. We first estimated the theoretical change in baseline DA elicited by the medication. For this, we
1290 computed the equivalent τ for the λ they report in the OFF and ON medication conditions
1291 ($\lambda_{\text{OFF}} = 1.51$, $\lambda_{\text{ON}} = 1.19$), using the relationship: $(1 - \tau) = \lambda$. We then computed the baseline
1292 DA levels that would give rise to the τ_{ON} and τ_{OFF} . With this, we computed the change in
1293 baseline DA (ΔDA) equivalent to the change $\Delta\tau = \tau_{\text{ON}} - \tau_{\text{OFF}}$. This ΔDA is the theoretical
1294 change in baseline DA elicited by the medication (Fig. 8F).
- 1295 2. To generate Fig. 8F, we sampled a set of λ from a Gaussian distribution centered at a mean of
1296 $\mu_{\lambda} = 1.51$ and a standard deviation of $\sigma_{\lambda}^2 = 3$, to emulate the distribution of λ_{OFF} they report in
1297 the OFF condition. We then computed the equivalent τ for that set of λ with the relationship
1298 above. We will call this the distribution of τ'_{OFF} .
- 1299 3. We used the derived τ'_{OFF} distribution to compute the equivalent dopamine levels. We imposed a
1300 change in baseline DA equal to the ΔDA computed in the first step and computed the new set of τ
1301 for that set of new baseline DA levels (τ'_{ON}). The ‘drug effect in loss aversion’ reported in Fig.
1302 8F is the $\tau'_{\text{ON}} - \tau'_{\text{OFF}}$ for each sample.

1303 5. Details on habenula lesion data

1304 5.1 Animals, surgery and lesions.

1305 The rodent data we re-analyzed here were first reported in Tian and Uchida (2015)⁵⁵. Below we provide a
1306 brief description of the methods. Further methodological details can be found in the original paper. 12
1307 mice were used. Bilateral habenula lesions were performed in five animals. Seven animals were in the
1308 control group including two with sham-lesion operation, one with only small contra-lateral side lesion of
1309 the medial habenula, and four animals without operations in the habenula. During surgery, a head plate
1310 was implanted on the skull, and adeno-associated virus (AAV) that express channelrhodopsin-2 (ChR2)
1311 in a Cre-dependent manner was injected into the VTA (from bregma: 3.1 mm posterior, 0.7 mm lateral,
1312 4–4.2 mm ventral). After recovery from surgery, mice were trained on the conditioning task, after which
1313 mice were randomly selected to be in lesion or sham-lesion group. Electrolytic lesions were made
1314 bilaterally using a stainless-steel electrode (15 kU, MicroProbes, MS301G) with a cathodal current of 150

1315 mA. Each side of the brain was lesioned at two locations (from bregma: 1.6 mm/1.9 mm posterior, 1.15
1316 mm lateral, 2.93 mm depth, with a 14 angle). For sham-lesion operations, no current was applied. In the
1317 same surgery, a microdrive containing electrodes and an optical fiber was implanted in the VTA (from
1318 bregma: 3.1 mm posterior, 0.7 mm lateral, 3.8–4.0 mm ventral)⁹⁸.

1319 **5.2 Behavioral task**

1320 Twelve mice were trained on a probabilistic Pavlovian task. Each trial the animal experienced one of four
1321 odor cues for 1 s, followed by a 1-s pause, followed by a reward (3.75 μ l water), an aversive air puff or
1322 nothing. Odor 1 to 3 signaled a 90%, 50% and 10% probability of reward, respectively. Odor 4 signaled a
1323 90% probability of air puff. Odor identities were randomized across trials and included: isoamyl acetate,
1324 eugenol, 1-hexanol, p-cymene, ethyl butyrate, 1-butanol, and carvone (1/10 dilution in paraffin oil). Inter-
1325 trial intervals were exponentially distributed. An infrared beam was positioned in front of the water
1326 delivery spout and each beam break was recorded as one lick event. We report the average lick rate over
1327 the interval 500–2,000 ms after cue onset.

1328 **5.3 Electrophysiology**

1329 Recordings were made using a custom-built microdrive equipped with 200- μ m-fiber optic-coupled with
1330 eight tetrodes. DA neurons were identified optogenetically⁹⁸. A stimulus-associated spike latency test
1331 (SALT) algorithm⁹⁹ was used to determine whether light pulses significantly changed a neuron's spike
1332 timing.

1333 **5.4 Neural data analysis**

1334 Data analyses were performed using MATLAB R2021b (Mathworks). To measure firing rates,
1335 peristimulus time histograms (PSTHs) were constructed using 1-ms bins. These histograms were then
1336 smoothed by convolving with the function $f(t) = (1 - e^{-t}) \cdot e^{-\frac{t}{\tau}}$ where τ was a time constant set to 20
1337 ms as in¹⁸. 44 dopamine neurons were recorded from lesioned animals (5 animals, 30 sessions), and 45
1338 dopamine neurons were recorded from control animals (7 animals, 35 sessions). We pooled all the cells
1339 across animals in each group for analysis. Cue-evoked responses were defined as the average activity
1340 from 0 to 400 ms after cue onset. Outcome-evoked responses were defined as the average activity from
1341 2000 to 2600 ms after cue onset.

1342 The normalization of cue response shown in Fig. 4 was carried out following a previous work³⁶ on a per-

1343 cell basis as: $c_{50}^{norm} = \frac{c_{50} - \bar{c}_{10}}{\bar{c}_{90} - \bar{c}_{10}}$, where \bar{c}_{90} , \bar{c}_{10} correspond to the mean across trials within a cell for the

1344 90% and 10% probability cure responses. To derive the t-statistics in Fig. 4d, we performed a two-tailed t-
1345 test of the cell's normalized responses to the 50% cue against the average midway point between
1346 responses to the 10% cue and responses to the 90% cue.

1347 The derivation of asymmetric scaling factors from outcome responses (τ_i), was carried out following³⁶,
1348 with some modifications to adapt it to the task. The procedure is illustrated in Extended Data Fig. 3.

1349 • To compute the reversal points, outcome responses were first aligned to the RPE for each trial
1350 type, computed with the true expected value of each reward distribution. Assuming a fixed reward
1351 value of 1 (arbitrary units), the expected value for the 90%, 50%, 10% reward probability trials
1352 corresponded to 0.1, 0.5, 0.9, respectively. Given this, omission responses from the 90%, 50%,
1353 10% reward probability trials correspond to RPEs of -0.9, -0.5 and -0.1. The rewarded responses
1354 from the 90%, 50%, 10% reward probability trials correspond to RPEs of 0.1, 0.5 and 0.9. The
1355 reward value is arbitrary and doesn't have an effect in this computation as it only shifts the RPE
1356 axis by a fixed amount. The reversal point for each cell (Z_i) was defined as the RPE that
1357 maximized the number of positive responses to RPEs greater than Z_i plus the number of negative
1358 responses to RPEs less than Z_i . The distribution of reversal points is reported in Extended Data
1359 Fig. 4. To obtain statistics for reliability of the computed reversal points, we partitioned the data
1360 into random halves and estimated the reversal point for each cell separately in each half. We
1361 repeated this procedure 1000 times with different random partitions, and we report the
1362 distribution of Pearson's correlation across these 1000 folds (Extended Data Fig. 4).

1363 • After measuring reversal points, we fit linear functions separately to the positive and negative
1364 domains. Given that dopamine's responses are non-linear in the reward space but present a
1365 putative utility function¹⁰⁰, we approximated the underlying utility function from the dopamine
1366 responses to RPEs of varying magnitudes. We used these empirical utilities instead of raw RPEs
1367 for computing the slopes that correspond to α_i^+ , α_i^- . We then computed the asymmetric scaling
1368 factors as $\tau_i = \frac{\alpha_i^+}{\alpha_i^+ + \alpha_i^-}$. We performed the same cross-validation procedure used for the reversal
1369 points. The distribution of R value across the 1000 folds are reported in Extended Data Fig. 4.

1370 A key prediction of distributional RL³⁶ is the presence of a correlation (across cells) between reversal
1371 points Z_i and asymmetric scaling factors τ_i . To elucidate whether signatures of distributional RL were
1372 still present after lesions, we followed the procedure given by Dabney et al. (2020)³⁶ to compute this
1373 correlation. We first randomly split the data into two disjoint halves of trials. In one half, we first

1374 calculated reversal points Z_i^1 and used them to calculate α_i^+, α_i^- . In the other half, we again calculated the
1375 reversal points Z_i^2 . The correlation we report in Extended Data Fig. 4 is between Z_i^2 and $\tau_i = \frac{\alpha_i^+}{\alpha_i^+ + \alpha_i^-}$.

1376 **5.5 Model fitting to the anticipatory licking responses**

1377 For each trial we computed the average lick rate over the interval 500–2,000 ms after cue onset. For each
1378 model, we fit the free parameters to the lick rates using maximum likelihood estimation. The optimization
1379 was performed using the SciPy optimization toolbox (Python) that minimized the difference between the
1380 predicted lick rates and the ground truth ones, with a uniform prior distribution over the parameters. The
1381 fits were done considering three RL models that had between 2 and 3 parameters. The models, parameters
1382 and bounds used for each of them are detailed in table 2.

1383

1384

1385 **Tables**

1386 **Table 1 – Biophysical model parameters**

1387

Parameter	Abbreviation	Value
DA association constant to D2 autorreceptors	k_{on}^{D2term}	$0.3 M^{-1}s^{-1}$
DA dissociation constant to D2 autorreceptors	k_{off}^{D2term}	$0.003 s^{-1}$
DA association constant to D1 receptors	k_{on}^{D1}	$0.01 M^{-1}s^{-1}$
DA dissociation constant to D1 receptors	k_{off}^{D1}	$10 s^{-1}$
DA association constant to D2 receptors	k_{on}^{D2}	$0.2 M^{-1}s^{-1}$
DA dissociation constant to D2 receptors	k_{off}^{D2}	$2 s^{-1}$
Release probability from terminals	P_r	$1 a. u.$
Release capacity from terminals	γ_{pr_n}	$2 a. u.$
D2 autorreceptor efficacy	α	$3 a. u.$
DAT maximal uptake capacity	$V_{max}^{pr_n}$	$1500 nMs^{-1}$
Michaelis-Menten parameter DAT-mediated DA uptake	K_m	$160 nM$
Constant for dopamine removal not mediated by DAT's	K_{nonDAT}	$0.04 nM$
Bromocriptine association constant to D2 autorreceptors	$k_{on}^{Drug,D2s}$	$0.02083 M^{-1}s^{-1}$
Bromocriptine dissociation constant to D2 autorreceptors	$k_{off}^{Drug,D2s}$	$0.1 s^{-1}$
Bromocriptine association constant to D2 receptors	$k_{on}^{Drug,D2l}$	$0.04 M^{-1}s^{-1}$
Bromocriptine dissociation constant to D2 receptors	$k_{off}^{Drug,D2l}$	$0.1 s^{-1}$

1388 a.u. = arbitrary units

1389 M = mols

1390 s = seconds

1391

1392

1393 Table 2- Reinforcement learning models fit to the behavioral data from Tian & Uchida

Model	Formulation	Parameters	Parameter bounds
TD learning	$\delta = r - V$ $V \leftarrow V + \alpha \cdot \delta$ $Licking = \beta \cdot V$	α, β	$\alpha \in [.001, 1]$ $\beta \in [.1, 10]$
TD learning with reward sensitivity	$\delta = \rho \cdot r - V$ $V \leftarrow V + \alpha \cdot \delta$ $Licking = \beta \cdot V$	α, ρ, β	$\alpha \in [.001, 1]$ $\rho \in [.001, 10]$ $\beta \in [.1, 10]$
Risk sensitive TD learning	$\delta = r - V$ $V \leftarrow V + \alpha^+ \cdot \delta \quad \text{if } \delta > 0$ $V \leftarrow V + \alpha^- \cdot \delta \quad \text{if } \delta < 0$ $Licking = \beta \cdot V$	$\alpha^+, \alpha^-, \beta$	$\alpha^+ \in [.001, 1]$ $\alpha^- \in [.001, 1]$ $\beta \in [.1, 10]$

1394

1395

1396 Data availability

1397 The neural data and simulation results reported in this article have been shared in a public
 1398 deposit source in: https://osf.io/cr5mv/?view_only=bd13a2d2de1947699b56ce70610b0e9b

1399

1400 Code availability

1401 The accession codes for the data as well as the code for analysis and simulations are available at:
 1402 https://github.com/sandraromerop/D1D2_Dopamine

1403

1404

1405 **References**

- 1406 1. Brown, V. M., Zhu L., Solway A., Wang M., McCurry K., King-Casas B. & Chiu P.
1407 Reinforcement Learning Disruptions in Individuals With Depression and Sensitivity to
1408 Symptom Change Following Cognitive Behavioral Therapy. *JAMA Psychiatry* **78**, 1113–
1409 1122 (2021).
- 1410 2. Groman, S. M., Thompson, S. L., Lee, D. & Taylor, J. R. Reinforcement learning detuned in
1411 addiction: integrative and translational approaches. *Trends Neurosci.* **45**, 96–105 (2022).
- 1412 3. Ligneul, R., Sescousse, G., Barbalat, G., Domenech, P. & Dreher, J.-C. Shifted risk
1413 preferences in pathological gambling. *Psychol. Med.* **43**, 1059–1068 (2013).
- 1414 4. Mason, L., O’Sullivan, N., Bentall, R. P. & El-Dereby, W. Better than I thought: positive
1415 evaluation bias in hypomania. *PLoS One* **7**, e47754 (2012).
- 1416 5. Pizzagalli, D. A., Iosifescu, D., Hallett, L. A., Ratner, K. G. & Fava, M. Reduced hedonic
1417 capacity in major depressive disorder: evidence from a probabilistic reward task. *J.*
1418 *Psychiatr. Res.* **43**, 76–87 (2008).
- 1419 6. Verdejo-Garcia, A., Chong, T. T.-J., Stout, J. C., Yücel, M. & London, E. D. Stages of
1420 dysfunctional decision-making in addiction. *Pharmacol. Biochem. Behav.* **164**, 99–105
1421 (2018).
- 1422 7. Lim, T. V., Cardinal, R. N., Bullmore, E. T., Robbins, T. W. & Ersche, K. D. Impaired
1423 Learning From Negative Feedback in Stimulant Use Disorder: Dopaminergic Modulation.
1424 *Int. J. Neuropsychopharmacol.* **24**, 867–878 (2021).
- 1425 8. Schönfelder, S., Langer, J., Schneider, E. E. & Wessa, M. Mania risk is characterized by an
1426 aberrant optimistic update bias for positive life events. *J. Affect. Disord.* **218**, 313–321
1427 (2017).
- 1428 9. Dayan, P. & Daw, N. D. Decision theory, reinforcement learning, and the brain. *Cogn. Affect.*
1429 *Behav. Neurosci.* **8**, 429–453 (2008).
- 1430 10. Katahira, K. The relation between reinforcement learning parameters and the influence of
1431 reinforcement history on choice behavior. *J. Math. Psychol.* **66**, 59–69 (2015).

- 1432 11. Sutton, R. S. & Barto, A. G. *Reinforcement Learning: An Introduction*. (Bradford Books,
1433 2018).
- 1434 12. Maia, T. V. & Frank, M. J. From reinforcement learning models to psychiatric and
1435 neurological disorders. *Nat. Neurosci.* **14**, 154–162 (2011).
- 1436 13. Eldar, E., Rutledge, R. B., Dolan, R. J. & Niv, Y. Mood as Representation of Momentum.
1437 *Trends Cogn. Sci.* **20**, 15–24 (2016).
- 1438 14. Rutledge, R. B., Skandali, N., Dayan, P. & Dolan, R. J. A computational and neural model of
1439 momentary subjective well-being. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 12252–12257 (2014).
- 1440 15. Floresco, S. B., West, A. R., Ash, B., Moore, H. & Grace, A. A. Afferent modulation of
1441 dopamine neuron firing differentially regulates tonic and phasic dopamine transmission. *Nat.*
1442 *Neurosci.* **6**, 968–973 (2003).
- 1443 16. Wang, Y., Toyoshima, O., Kunitatsu, J., Yamada, H. & Matsumoto, M. Tonic firing mode
1444 of midbrain dopamine neurons continuously tracks reward values changing moment-by-
1445 moment. *Elife* **10**, (2021).
- 1446 17. Schultz, W., Dayan, P. & Montague, P. R. A neural substrate of prediction and reward.
1447 *Science* **275**, 1593–1599 (1997).
- 1448 18. Eshel, N., Tian, J., Bukwich, M. & Uchida, N. Dopamine neurons share common response
1449 function for reward prediction error. *Nat. Neurosci.* **19**, 479–486 (2016).
- 1450 19. Steinberg, E. Keiflin, R., Boivin, J., Witten I., Deisseroth K. & Janak P. A causal link
1451 between prediction errors , dopamine neurons and learning. *Nat. Neurosci.* **16**, 966–973
1452 (2013).
- 1453 20. Waelti, P., Dickinson, A. & Schultz, W. Dopamine responses comply with basic assumptions
1454 of formal learning theory. *Nature* **412**, 43–48 (2001).
- 1455 21. Korn, C. W., Sharot, T., Walter, H., Heekeren, H. R. & Dolan, R. J. Depression is related to
1456 an absence of optimistically biased belief updating about future life events. *Psychol. Med.* **44**,
1457 579–592 (2014).
- 1458 22. Rutledge, R. B., Lazzaro, S. C., Lau, B., Myers, C. E., Gluck, M. A., & Glimcher, P. W.
1459 Dopaminergic drugs modulate learning rates and perseveration in Parkinson’s patients in a
1460 dynamic foraging task. *J. Neurosci.* **29**, 15104–15114 (2009).

- 1461 23. Frank, M. J., Seeberger, L. C. & O'Reilly, R. C. By carrot or by stick: Cognitive
1462 reinforcement learning in Parkinsonism. *Science* **306**, 1940–1943 (2004).
- 1463 24. Cools, R., Frank M.J., Gibbs S., Miyakawa A., Jagust W. & D'Esposito M. Striatal
1464 dopamine predicts outcome-specific reversal learning and its sensitivity to dopaminergic
1465 drug administration. *J. Neurosci.* **29**, 1538–1543 (2009).
- 1466 25. Timmer, M. H. M., Sescousse, G., van der Schaaf, M. E., Esselink, R. A. J. & Cools, R.
1467 Reward learning deficits in Parkinson's disease depend on depression. *Psychol. Med.* **47**,
1468 2302–2311 (2017).
- 1469 26. Gradin, V. B., Kumar, P., Waiter, G., Ahearn, T., Stickle, C., Milders, M., Reid, I., Hall, J.,
1470 & Steele, J. D. Expected value and prediction error abnormalities in depression and
1471 schizophrenia. *Brain* **134**, 1751–1764 (2011).
- 1472 27. Kumar, P., Waiter, G., Ahearn, T., Milders, M., Reid, I., & Steele, J. D. Abnormal temporal
1473 difference reward-learning signals in major depression. *Brain* **131**, 2084–2093 (2008).
- 1474 28. Pizzagalli, D. A., Holmes, A. J., Dillon, D. G., Goetz, E. L., Birk, J. L., Bogdan, R.,
1475 Dougherty, D. D., Iosifescu, D. V., Rauch, S. L., & Fava, M. Reduced caudate and nucleus
1476 accumbens response to rewards in unmedicated individuals with major depressive disorder.
1477 *Am. J. Psychiatry* **166**, 702–710 (2009).
- 1478 29. Robinson, O. J., Cools, R., Carlisi, C. O., Sahakian, B. J. & Drevets, W. C. Ventral striatum
1479 response during reward and punishment reversal learning in unmedicated major depressive
1480 disorder. *Am. J. Psychiatry* **169**, 152–159 (2012).
- 1481 30. Collins, A. G. E. & Frank, M. J. Opponent actor learning (OpAL): modeling interactive
1482 effects of striatal dopamine on reinforcement learning and choice incentive. *Psychol. Rev.*
1483 **121**, 337–366 (2014).
- 1484 31. Frank, M. J. Dynamic Dopamine Modulation in the Basal Ganglia: A Neurocomputational
1485 Account of Cognitive Decits in Medicated and Non-medicated Parkinsonism. *J. Cogn.*
1486 *Neuroci.* **17**, 51-72 (2005).
- 1487 32. Mikhael, J. G. & Bogacz, R. Learning Reward Uncertainty in the Basal Ganglia. *PLoS*
1488 *Comput. Biol.* **12**, 1–28 (2016).

- 1489 33. Yagishita, S., Hayashi- Takagi, A., Ellis-Davies, G.C.R., Urakubo, H., Ishii, S. & Kasai, H.
1490 A critical time window for dopamine actions on the structural plasticity of dendritic spines.
1491 *Science* **345**, 1616–1620 (2014).
- 1492 34. Iino, Y. Sawada, T., Yamaguchi, K. *et al.* Dopamine D2 receptors in discrimination learning
1493 and spine enlargement. *Nature* **579**, 555–560 (2020).
- 1494 35. Lee, S. J., Lodder, B., Chen, Y., Patriarchi, T., Tian, L & Sabatini, B. Cell-type-specific
1495 asynchronous modulation of PKA by dopamine in learning. *Nature* **590**, 451–456 (2021).
- 1496 36. Dabney, W., Kurth-Nelson, Z., Uchida, N., Starkweather, C., Hassabis, D., Munos, R. &
1497 Botvinick, M. A distributional code for value in dopamine- based reinforcement learning.
1498 *Nature* **1**, (2019).
- 1499 37. Bellemare, M. G., Dabney, W. & Munos, R. A distributional perspective on reinforcement
1500 learning. *34th International Conference on Machine Learning, ICML 2017* **1**, 693–711
1501 (2017).
- 1502 38. Lowet, A. S., Zheng, Q., Matias, S., Drugowitsch, J. & Uchida, N. Distributional
1503 Reinforcement Learning in the Brain. *Trends Neurosci.* **43**, 980–997 (2020).
- 1504 39. Rowland, M., Dadashi, R., Kumar, S., Munos, R., Bellemare, M. G. & Dabney, W. Statistics
1505 and samples in distributional reinforcement learning. *36th International Conference on*
1506 *Machine Learning, ICML 2019* **6**, 9727–9750 (2019).
- 1507 40. Mihatsch, O. & Neuneier, R. Risk-sensitive reinforcement learning. *Mach. Learn.* **49**, 267–
1508 290 (2002).
- 1509 41. Bellemare, M. G. & Dabney, W. *Distributional reinforcement learning*. (MIT Press, 2023).
- 1510 42. Jones, M. C. Expectiles and M-quantiles are quantiles. *Stat. Probab. Lett.* **20**, 149–153
1511 (1994).
- 1512 43. Houk, J., Davis, J., & Beiser, D. *Models of information processing in the basal ganglia*.
1513 (Bradford Books, 2019).
- 1514 44. Gerfen, C. The neostriatal mosaic: Multiple levels of compartmental organization in the basal
1515 ganglia. *Annu. Rev. Neurosci.* **15**, 285–320 (1992).

- 1516 45. Smith, Y., Bevan, M. D., Shink, E. & Bolam, J. P. Microcircuitry of the direct and indirect
1517 pathways of the basal ganglia. *Neuroscience* **86**, 353–387 (1998).
- 1518 46. Kravitz, A. V., Tye, L. D. & Kreitzer, A. C. Distinct roles for direct and indirect pathway
1519 striatal neurons in reinforcement. *Nat. Neurosci.* **15**, 816–818 (2012).
- 1520 47. Reynolds, J. N. J., Hyland, B. I. & Wickens, J. R. A cellular mechanism of reward-related
1521 learning. *Nature* **413**, 67–70 (2001).
- 1522 48. Gerfen, C. R. & Surmeier, D. J. Modulation of striatal projection systems by dopamine.
1523 *Annu. Rev. Neurosci.* **34**, 441–466 (2011).
- 1524 49. Richfield, E. K., Penney, J. B. & Young, A. B. Anatomical and affinity state comparisons
1525 between dopamine D1 and D2 receptors in the rat central nervous system. *Neuroscience* **30**,
1526 767–777 (1989).
- 1527 50. Rice, M. E. & Cragg, S. J. Dopamine spillover after quantal release: Rethinking dopamine
1528 transmission in the nigrostriatal pathway. *Brain Res. Rev.* **58**, 303–313 (2008).
- 1529 51. Gonon, F. G. & Buda, M. J. Regulation of dopamine release by impulse flow and by
1530 autoreceptors as studied by in vivo voltammetry in the rat striatum. *Neuroscience* **14**, 765–
1531 774 (1985).
- 1532 52. Dodson, P. D., Dreyer, J. K., Jennings, K. A., Syed, E. C. J., Wade-Martins, R., Cragg, S. J.,
1533 Bolam, J. P. & Magill, P. J. Representation of spontaneous movement by dopaminergic
1534 neurons is cell-type selective and disrupted in parkinsonism. *Proc. Natl. Acad. Sci. U. S. A.*
1535 **113**, E2180–E2188 (2016).
- 1536 53. Marcott, P. F., Mamaligas, A. A. & Ford, C. P. Phasic dopamine release drives rapid
1537 activation of striatal D2-receptors. *Neuron* **84**, 164–176 (2014).
- 1538 54. Jaskir, A. & Frank, M. J. On the normative advantages of dopamine and striatal opponency
1539 for learning and choice. *Elife* (2023).
- 1540 55. Tian, J. & Uchida, N. Habenula Lesions Reveal that Multiple Mechanisms Underlie
1541 Dopamine Prediction Errors. *Neuron* **87**, 1304–1316 (2015).
- 1542 56. Cui, Y., Yang, Y., Ni, Z., Dong, Y., Cai, G., Foncelle, A., Ma, S., Sang, K., Tang, S., Li, Y.,
1543 Shen, Y., Berry, H., Wu, S. & Hu, H. Astroglial Kir4.1 in the lateral habenula drives
1544 neuronal bursts in depression. *Nature* **554**, 323–327 (2018).

- 1545 57. Li, B., Piriz, J., Mirrione, M., Chung, C., D. Proulx, C., Schulz, D., Henn, F. & Malinow, R.
1546 Synaptic potentiation onto habenula neurons in the learned helplessness model of depression.
1547 *Nature* **470**, 535–541 (2011).
- 1548 58. Yang, Y., Cui, Y., Sang, K., Dong, Y., Ni, Z., Ma, S. & Hu, H. Ketamine blocks bursting in
1549 the lateral habenula to rapidly relieve depression. *Nature* **554**, 317–322 (2018).
- 1550 59. Dreyer, J. K., Herrik, K. F., Berg, R. W. & Hounsgaard, J. D. Influence of phasic and tonic
1551 dopamine release on receptor activation. *J. of Neurosci.* **30**, 14273–14283 (2010).
- 1552 60. Mikhael, J. G., Lai, L. & Gershman, S. J. *Rational inattention and tonic dopamine. PLoS*
1553 *Comput. Biol.* **17** e1008659 (2021).
- 1554 61. Vingerhoets, F. J., Snow, B. J., Schulzer, M., Morrison, S., Ruth, T. J., Holden, J. E., Cooper,
1555 S. & Calne, D. B. Reproducibility of fluorine-18-6-fluorodopa positron emission tomography
1556 in normal human subjects. *J. Nucl. Med.* **35**, 18–24 (1994).
- 1557 62. Cools, R., Gibbs, S. E., Miyakawa, A., Jagust, W. & D’Esposito, M. Working memory
1558 capacity predicts dopamine synthesis capacity in the human striatum. *J. Neurosci.* **28**, 1208–
1559 1212 (2008).
- 1560 63. Hoffmann, I. S. & Cubeddu, L. X. Differential effects of bromocriptine on dopamine and
1561 acetylcholine release modulatory receptors. *J. Neurochem.* **42**, 278–282 (1984).
- 1562 64. Tissari, A. H., Rossetti, Z. L., Meloni, M., Frau, M. I. & Gessa, G. L. Autoreceptors mediate
1563 the inhibition of dopamine synthesis by bromocriptine and lisuride in rats. *Eur. J.*
1564 *Pharmacol.* **91**, 463–468 (1983).
- 1565 65. Lieberman, A. Depression in Parkinson’s disease -- a review. *Acta Neurol. Scand.* **113**, 1–8
1566 (2006).
- 1567 66. Leentjens, A. F. G., Van den Akker, M., Metsemakers, J. F. M., Lousberg, R. & Verhey, F.
1568 R. J. Higher incidence of depression preceding the onset of Parkinson’s disease: a register
1569 study. *Mov. Disord.* **18**, 414–418 (2003).
- 1570 67. Nilsson, F. M., Kessing, L. V., Sørensen, T. M., Andersen, P. K. & Bolwig, T. G. Major
1571 depressive disorder in Parkinson’s disease: a register-based study. *Acta Psychiatr. Scand.*
1572 **106**, 202–211 (2002).

- 1573 68. Remy, P., Doder, M., Lees, A., Turjanski, N. & Brooks, D. Depression in Parkinson's
1574 disease: loss of dopamine and noradrenaline innervation in the limbic system. *Brain* **128**,
1575 1314–1322 (2005).
- 1576 69. Weintraub, D., Newberg, A. B., Cary, M. S., Siderowf, A. D., Moberg, P. J., Kleiner-Fisman,
1577 G., Duda, J. E., Stern, M. B., Mozley, D. & Katz, I. R. Striatal dopamine transporter imaging
1578 correlates with anxiety and depression symptoms in Parkinson's disease. *J. Nucl. Med.* **46**,
1579 227–232 (2005).
- 1580 70. Kish, S. J., Shannak, K. & Hornykiewicz, O. Uneven pattern of dopamine loss in the striatum
1581 of patients with idiopathic Parkinson's disease. Pathophysiologic and clinical implications.
1582 *N. Engl. J. Med.* **318**, 876–880 (1988).
- 1583 71. Timmer, M. H. M., Sescousse, G., Esselink, R. A. J., Piray, P. & Cools, R. Mechanisms
1584 Underlying Dopamine-Induced Risky Choice in Parkinson's Disease With and Without
1585 Depression (History). *Comput Psychiatr* **2**, 11–27 (2018).
- 1586 72. Hikida, T., Kimura, K., Wada, N., Funabiki, K. & Nakanishi Shigetada, S. Distinct Roles of
1587 Synaptic Transmission in Direct and Indirect Striatal Pathways to Reward and Aversive
1588 Behavior. *Neuron* **66**, 896–907 (2010).
- 1589 73. Hikida, T., Yawata, S., Yamaguchi, T., Danjo, T., Sasaoka, T., Wang, Y. & Nakanishi S.
1590 Pathway-specific modulation of nucleus accumbens in reward and aversive behavior via
1591 selective transmitter receptors. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 342–347 (2013).
- 1592 74. Danjo, T., Yoshimi, K., Funabiki, K., Yawata, S. & Nakanishi, S. Aversive behavior induced
1593 by optogenetic inactivation of ventral tegmental area dopamine neurons is mediated by
1594 dopamine D2 receptors in the nucleus accumbens. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 6455–
1595 6460 (2014).
- 1596 75. Yamaguchi, T., Goto, A., Nakahara, I., Yawata, S., Hikida, T., Matsuda, M., Funabiki, K. &
1597 Nakanishi, S. Role of PKA signaling in D2 receptor-expressing neurons in the core of the
1598 nucleus accumbens in aversive learning. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 11383–11388
1599 (2015).
- 1600 76. Bayer, H. M. & Glimcher, P. W. Midbrain dopamine neurons encode a quantitative reward
1601 prediction error signal. *Neuron* **47**, 129–141 (2005).

- 1602 77. Hart, A. S., Rutledge, R. B., Glimcher, P. W. & Phillips, P. E. M. Phasic dopamine release in
1603 the rat nucleus accumbens symmetrically encodes a reward prediction error term. *J.*
1604 *Neurosci.* **34**, 698–704 (2014).
- 1605 78. Grace, A. A. Dysregulation of the dopamine system in the pathophysiology of schizophrenia
1606 and depression. *Nat. Rev. Neurosci.* **17**, 524–532 (2016).
- 1607 79. Anstrom, K. K., Miczek, K. A. & Budygin, E. A. Increased phasic dopamine signaling in the
1608 mesolimbic pathway during social defeat in rats. *Neuroscience* **161**, 3–12 (2009).
- 1609 80. Razzoli, M., Andreoli, M., Michielin, F., Quarta, D. & Sokal, D. M. Increased phasic activity
1610 of VTA dopamine neurons in mice 3 weeks after repeated social defeat. *Behav. Brain Res.*
1611 **218**, 253–257 (2011).
- 1612 81. Markovic, T., Pederson, C. E., Massaly, N., Vachez, Y., Ruyle, B., Murphy, C. A.,
1613 Abiraman, K., Hoon Shin, J., Garcia, J. H., Jean Yoon, H., Alvarez, V. A., Bruchas, M. R.,
1614 Creed, M. C. & Moron, J. A. Pain induces adaptations in ventral tegmental area dopamine
1615 neurons to drive anhedonia-like behavior. *Nat. Neurosci.* **24**, 1601–1613 (2021).
- 1616 82. Guo, Z., Li, S., Wu, J., Zhu, X. & Zhang, Y. Maternal deprivation increased vulnerability to
1617 depression in adult rats through DRD2 promoter methylation in the ventral tegmental area.
1618 *Front. Psychiatry* **13**, 827667 (2022).
- 1619 83. Peng, B., Hu, Q., Liu, J., Guo, S., Borgland, S. L. & Liu, S. Corticosterone attenuates
1620 reward-seeking behavior and increases anxiety via D2 receptor signaling in ventral tegmental
1621 area dopamine neurons. *J. Neurosci.* **41**, 1566–1581 (2021).
- 1622 84. Tye, K. M., Mirzabekov, J. J., Warden, M. R., Ferenczi, E. A., Tsai, H-C., Finkelstein, J.,
1623 Kim, S-Y., Adhikari, A., Thompson, K. R., Andalman A. S., Gunaydin, L. A., Witten I. &
1624 Deisseroth K. Dopamine neurons modulate neural encoding and expression of depression-
1625 related behaviour. *Nature* **493**, 537–541 (2013).
- 1626 85. Baek, K., Kwon, J., Chae, J. H., Chung, Y. A., Kralik, J. D., Min, J. A., Huh, H., Choi, K.
1627 M., Jang, K. I., Lee, N. B., Kim, S., Peterson, B. S., & Jeong, J. Heightened aversion to risk
1628 and loss in depressed patients with a suicide attempt history. *Sci. Rep.* **7**, 11228 (2017).

- 1629 86. Smoski, M. J., Lynch, T. R., Rosenthal, M. Z., Cheavens, J. S., Chapman, A. L., & Krishnan,
1630 R. R. Decision-making and risk aversion among depressive adults. *J. Behav. Ther. Exp.*
1631 *Psychiatry* **39**, 567–576 (2008).
- 1632 87. van Holst, R. J., Sescousse, G., Janssen, L. K., Janssen, M., Berry, A. S., Jagust, W. J., &
1633 Cools, R. Increased Striatal Dopamine Synthesis Capacity in Gambling Addiction. *Biol.*
1634 *Psychiatry* **83**, 1036–1043 (2018).
- 1635 88. Cools, R., Altamirano, L. & D’Esposito, M. Reversal learning in Parkinson’s disease
1636 depends on medication status and outcome valence. *Neuropsychologia* **44**, 1663–1673
1637 (2006).
- 1638 89. Cools, R., Barker, R. A., Sahakian, B. J. & Robbins, T. W. Enhanced or impaired cognitive
1639 function in Parkinson’s disease as a function of dopaminergic medication and task demands.
1640 *Cereb. Cortex* **11**, 1136–1143 (2001).
- 1641 90. Cools, R., Barker, R. A., Sahakian, B. J. & Robbins, T. W. L-Dopa medication remediates
1642 cognitive inflexibility, but increases impulsivity in patients with Parkinson’s disease.
1643 *Neuropsychologia* **41**, 1431–1441 (2003).
- 1644 91. Swainson, R., Rogers, R. D., Sahakian, B. J., Summers, B. A., Polkey, C. E., & Robbins, T.
1645 W. Probabilistic learning and reversal deficits in patients with Parkinson’s disease or frontal
1646 or temporal lobe lesions: possible adverse effects of dopaminergic medication.
1647 *Neuropsychologia* **38**, 596–612 (2000).
- 1648 92. Cox, S. M., Frank, M. J., Larcher, K., Fellows, L. K., Clark, C. A., Leyton, M., & Dagher, A.
1649 Striatal D1 and D2 signaling differentially predict learning from positive and negative
1650 outcomes. *Neuroimage* **109**, 95–101 (2015).
- 1651 93. Savitz, J. B. & Drevets, W. C. Neuroreceptor imaging in depression. *Neurobiol. Dis.* **52**, 49–
1652 65 (2013).
- 1653 94. Rescorla, R. A. & Wagner, A. R. *A theory of Pavlovian conditioning: Variations in the*
1654 *effectiveness of reinforcement and nonreinforcement, Classical Conditioning II.* 64–99
1655 (1972).

- 1656 95. Kim, H. R., Malik, A. N., Mikhael, J. G., Bech, P., Tsutsui-Kimura, I., Sun, F., Zhang, Y., Li,
1657 Y., Watabe-Uchida, M., Gershman, S. J. & Uchida, N. A unified framework for dopamine
1658 signals across timescales. *Cell* **183**, 1600-1616.e25 (2020).
- 1659 96. Bertsekas, D. P. & Tsitsiklis, J. *Neuro-Dynamic Programming*. (Athena Scientific, 1996).
- 1660 97. Mierau, J., Schneider, F. J., Ensinger, H. A., Chio, C. L., Lajiness, M. E., & Huff, R. M.
1661 Pramipexole binding and activation of cloned and expressed dopamine D2, D3 and D4
1662 receptors. *Eur. J of Pharmac.* **290**, 29-36 (1995).
- 1663 98. Cohen, J. Y., Haesler, S., Vong, L., Lowell, B. B. & Uchida, N. Neuron-type-specific signals
1664 for reward and punishment in the ventral tegmental area. *Nature* **482**, 85–88 (2012).
- 1665 99. Kvitsiani, D., Ranade, S., Hangya, B., Taniguchi, H., Huang, J. Z. & Kepec, A. Distinct
1666 behavioural and network correlates of two interneuron types in prefrontal cortex. *Nature* **498**,
1667 363–366 (2013).
- 1668 100. Stauffer, W. R., Lak, A. & Schultz, W. Dopamine reward prediction error responses reflect
1669 marginal utility. *Curr. Biol.* **24**, 2491–2500 (2014).
- 1670
- 1671
- 1672
- 1673
- 1674

1675 **Acknowledgements**

1676 We thank Ju Tian for providing the data used in this study; Jakob Dreyer for making the
1677 biophysical model available for open-source use; Roshan Cools, Peter Dayan and Sam Gershman
1678 for valuable feedback on the manuscript; Dan Polley, Mark Andermann and Elizabeth Phelps for
1679 mentorship and helpful discussions. We thank Bernardo Sabatini for discussions, and Isobel
1680 Green, Paul Masset, Mitsuko Watabe-Uchida and all the rest of the Uchida lab members for
1681 helpful comments at the different stages of this work. This work is supported by NIH BRAIN
1682 Initiative grants (R01NS116753, U19NS113201), the Simons Collaboration on Global Brain,
1683 and the Bipolar Disorder Seed Grant from the Harvard Brain Science Institute.

1684

1685 **Author information**

1686 **Contributions**

1687 S.R.P. and N.U. conceived the project. S.R.P performed the modeling work. S.R.P. wrote the
1688 first draft and S.R.P and N.U. edited the paper.

1689

1690 **Corresponding author**

1691 Correspondence to Naoshige Uchida (uchida@mcb.harvard.edu) and Sandra Romero Pinto
1692 (sromeropinto@g.harvard.edu).

1693

1694

1695

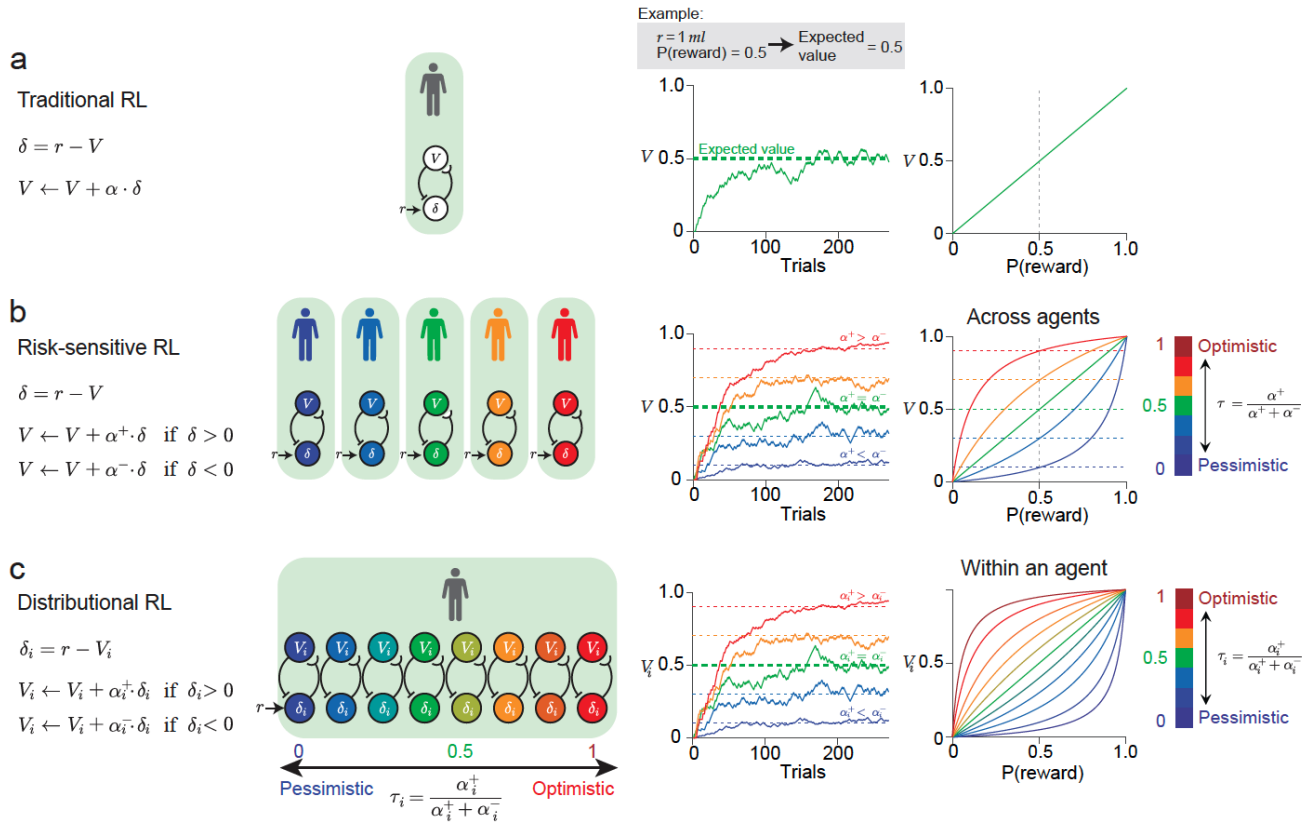
1696 **Competing interest statement**

1697 The authors declare no competing interests.

1698

1699

1700 **Figures**



1701

1702 **Figure 1 | Reinforcement learning models.**

1703 **a.** Traditional reinforcement learning with a single learning rate (α) for both positive and
 1704 negative RPEs (δ) for the value updates (left). This update rule makes value estimate (V)
 1705 converge on the expected value of the reward distribution (middle). When the reward probability
 1706 is varied (i.e., for Bernoulli distributions), the V at convergence scales linearly with the reward
 1707 probability (right).

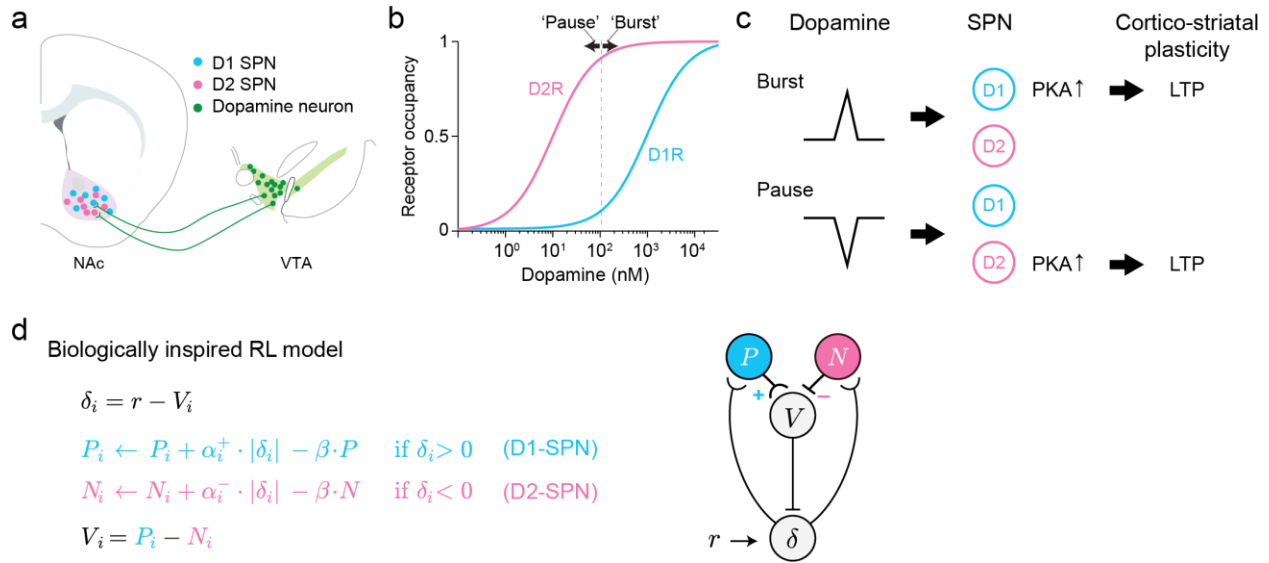
1708 **b.** Risk-sensitive reinforcement learning with different learning rates (α^+ , α^-) for positive and
 1709 negative RPEs, respectively (left). This update rule makes value estimate (V) converge on the
 1710 quantities that are higher or lower than the expected value of the reward distribution (middle). As
 1711 the reward probabilities are varied, the convexity of the convergent value V changes depending
 1712 on the asymmetry between α^+ and α^- (Methods 1.3.3). The level of the bias is determined by
 1713 the asymmetric learning rate parameter τ (right).

1714 **c.** Distributional reinforcement learning contains a set of value predictors (V_i) each with a given
 1715 learning rate for positive and negative RPEs (α_i^+ , α_i^- , respectively) (left). This makes each value

1716 predictor converge on the quantity equal to the τ_i -th expectile of the reward distribution. Thus,
1717 each value V_i represents an expectile, and together the set of V_i represents the entire distribution
1718 (Methods 1.2) (right).

1719

1720



1721

1722 **Figure 2 | Biologically inspired reinforcement learning model.**

1723 **a.** Schematic of the basal ganglia circuitry. Dopaminergic neurons in the VTA modulate
1724 plasticity at the level of the cortico-striatal synapses on SPNs in the NAc. The SPNs are
1725 subdivided depending on the dopamine receptor type they express (D1R or D2R).

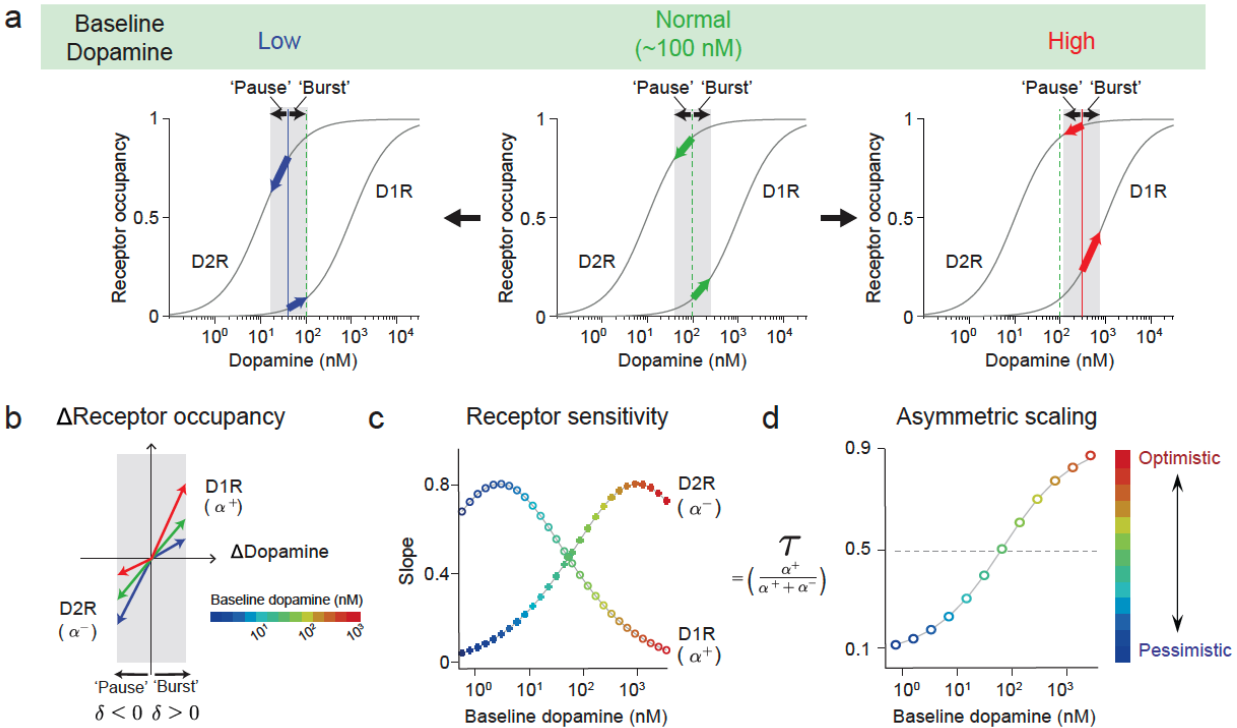
1726 **b.** Dose-occupancy curves for the D1R and D2R describing receptor occupancies as a function of
1727 dopamine concentrations. The curves are shifted between each other due to the different
1728 affinities of the receptors. The arrows represent 3-fold increase (“burst”) and decrease (“pause”)
1729 in dopamine concentrations, which causes left-ward or right-ward shifts of the same magnitudes
1730 in the log-scale.

1731 **c.** Schematic of the plasticity rules of VTA-NAc circuitry^{33–35}. Transient increases in dopamine,
1732 caused by bursts in firing rate of dopamine neurons, generates increases in PKA activity in D1R-
1733 expressing SPNs, leading to LTP in the cortico-striatal synapses. Transient decreases in
1734 dopamine, caused by pauses in firing rate of dopamine neurons, generates increases in PKA
1735 activity in D2R-expressing SPNs, leading to LTP in the cortico-striatal synapses.

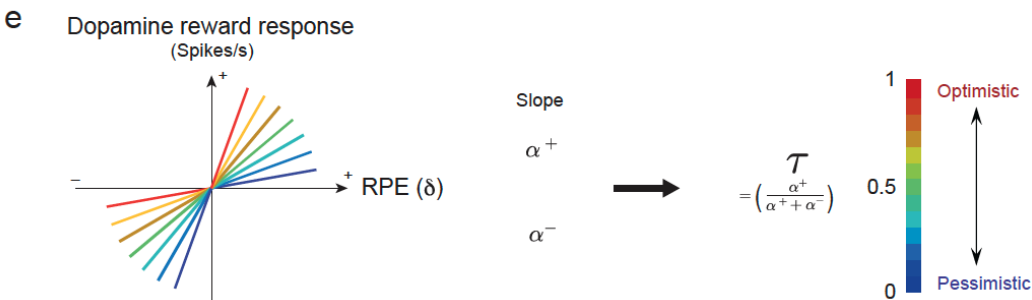
1736 **d.** Schematic and equations of biologically inspired reinforcement learning model³²

1737 VTA, ventral tegmental area; NAc, nucleus accumbens; SPN, spiny projection neurons; D1R,
1738 D1-type dopamine receptor; D2R, D2-type dopamine receptor; PKA, protein kinase A; LTP,
1739 long-term potentiation.

Model 1 Baseline dopamine → Asymmetric learning rates (τ)



Model 2 Different scaling of phasic dopamine → Asymmetric learning rates (τ)



1740

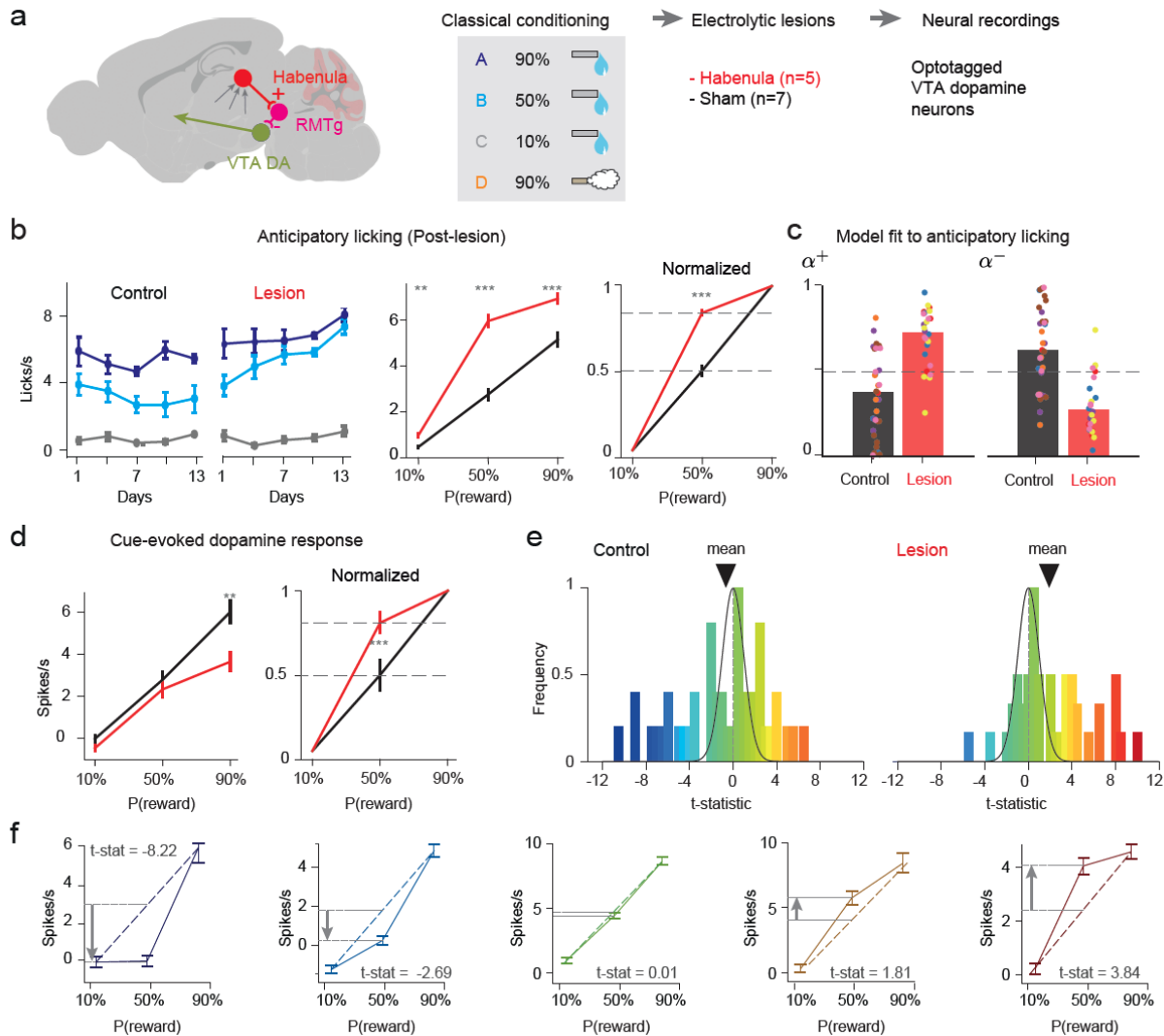
1741 **Figure 3 | Potential mechanisms for asymmetric learning.**

1742 **a.** Schematic of the mechanism by which increases or decreases in baseline dopamine modulates
 1743 the degree to which bursts and pauses in dopamine causes changes in D1R and D2R occupancy.

1744 Increases in baseline dopamine makes dopamine pauses to cause greater decreases in D2R
 1745 occupancy than the increases in D1R occupancy caused by dopamine bursts. Conversely,

1746 decreases in dopamine, makes dopamine bursts to cause smaller increases in D1R occupancy
 1747 than the decreases in D2R occupancy caused by dopamine pauses.

- 1748 **b.** Schematic of the change in receptor occupancies in D1R and D2R, for a given transient
1749 increase (‘burst’) or decrease (‘pause’) in dopamine, receptively. A pause and a burst in
1750 dopamine correspond to $\delta < 0$ and $\delta > 0$ in the model. The slope is modulated by the baseline
1751 dopamine (colormap) and corresponds to the receptor’s sensitivity to dopamine transients.
- 1752 **c.** Receptor sensitivity for D1R and D2R as a function of baseline dopamine. In Model 1, we
1753 assume that the receptor sensitivity acts as a scaling factor on the PKA activity induced by burst
1754 and pauses. That is, $\text{PKA}_{\text{D1}} \propto \alpha^+ \cdot \delta \cdot \mathbf{1}_{\delta > 0}$ and $\text{PKA}_{\text{D2}} \propto \alpha^- \cdot \delta \cdot \mathbf{1}_{\delta < 0}$.
- 1755 **d.** Asymmetric scaling factor (τ) as a function of baseline dopamine. Colors depict how
1756 ‘optimistic’ or ‘pessimistic’ the convergent value estimate will be when learning with a given τ .
- 1757 **e.** Model 2. Left, the relationship between dopamine reward responses (spikes/s) and RPEs. The
1758 slopes of these response functions correspond to the asymmetric learning rates (α^+ , α^-) for
1759 positive and negative RPEs, respectively. Colors depict how ‘optimistic’ or ‘pessimistic’ the
1760 convergent value estimate will be when learning with a given asymmetric scaling factor.
- 1761
- 1762



1763

1764 **Figure 4 | Habenula lesions leads to optimistic reward-seeking behavior and cue-evoked**
 1765 **responses in dopamine neurons.**

1766 **a.** Schematic of the experiment performed by Tian and Uchida (2015)⁵⁵. Animals were trained in
 1767 a classical conditioning task in which 3 odor cues predicted rewards of different probabilities
 1768 (10%, 50%, 90%) and one odor cue predicted 80% probability of an air puff. Animals then
 1769 underwent habenula ($n = 5$) or sham ($n = 7$) lesions and trained on the task again. The neural
 1770 recordings were performed from optotagged VTA dopamine neurons once behavior stabilized.

1771 **b.** Anticipatory licking across sessions after lesions (left.). There was a significant increase in
 1772 anticipatory licking to the 10% (U-statistic = -2.895 , $P = 0.003792$, two-sided Mann-Whitney U-
 1773 test), 50% (U-statistic = -5.579 , $P < 1 \times 10^{-9}$, two-sided Mann-Whitney U-test) and 90% (U-

1774 statistic = -3.487 , $P = 0.00048$, two-sided Mann-Whitney U-test) cues ($n = 31$ for control $n = 30$
1775 for lesion) that results from progressive changes across sessions. The anticipatory licking curves
1776 show a linear scaling with reward probability in the control group, and a convex curve for the
1777 lesion group (mean \pm s.e.m across animals, U-statistic = -6.444 , $P < 1 \times 10^{-7}$, two-sided Mann-
1778 Whitney U-test for the 50% cue normalized response). These curves are predicted by RL agents
1779 with symmetric and asymmetric ($\alpha^+ > \alpha^-$) learning rates for the control and lesion groups,
1780 respectively, assuming a linear mapping between anticipatory licking and value prediction.

1781 **c.** RL model fits to the anticipatory licking on a trial-by-trial basis using a risk-sensitive RL
1782 models that allows for separate learning rates of positive and negative RPEs. Each dot represents
1783 a session ($n = 35$ control, $n = 30$ lesion) and each color a mouse ($n = 7$ control, $n = 5$ lesion). The
1784 fits show a significant difference in the learning rates between control and lesion groups (U-
1785 statistic = -4.679 , $P < 1.0 \times 10^{-5}$, pooling sessions across mice in each group).

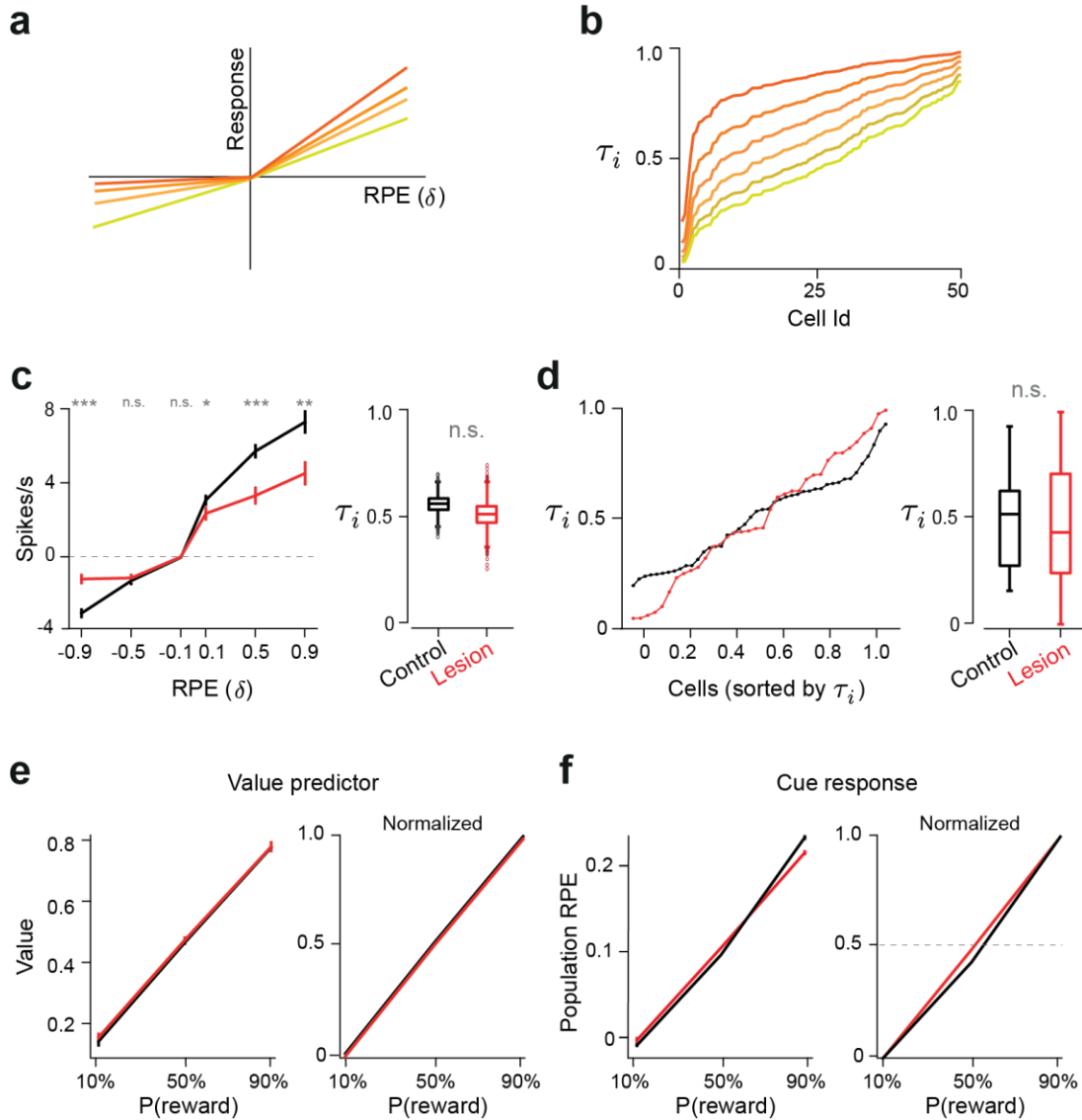
1786 **d.** Cue-evoked dopamine responses from opto-tagged VTA dopamine neurons (mean \pm s.e.m
1787 across neurons, $n = 45$ control group, $n = 44$ lesion group). There was a decrease in the absolute
1788 magnitude of responses to the 90% cue (U-statistic = 3.249 , $P = 0.0011$, two-sided Mann-
1789 Whitney U-test) after habenula lesions (left). The normalized cue-evoked responses show the
1790 similar pattern as the normalized anticipatory-licking with a linear and convex function for the
1791 control and lesion groups, respectively, with a significant increase in normalized response to the
1792 50% cue after lesions (U-statistic = -3.824 , $P = 0.000131$, two-sided Mann-Whitney U-test)
1793 These curves are predicted by agents with symmetric and asymmetric learning rates for control
1794 and lesion groups, respectively.

1795 **e.** Distribution of t-statistics comparing the cue-evoked response to the linear interpolation point
1796 between the 90% and 10% cue-evoked responses for each dopamine neuron. The distribution of
1797 t-statistics for the control and lesion cases was wider than what is expected from random noise
1798 (Monte Carlo test for standard deviation different from zero: $P = 0.0222$ control, $P = 0.0217$
1799 lesion, 1000 batches). The distribution was shifted to values larger than 0 in the lesion case
1800 (Monte Carlo test for mean larger than zero: $P = 1$ control, $P = 0.022$ lesion, 1000 batches)
1801 indicative of an optimistic bias in the distribution. The lesion group distribution was also
1802 significantly shifted to higher values with respect to the control group distribution (U-statistic =
1803 -2.815 , $P = 0.0024$, single-sided Mann-Whitney U-test). Arrow heads: the mean of the t -
1804 statistics.

1805 **f.** Example of t -statistics calculations for dopamine neurons taken from the control group (mean
1806 \pm s.e.m across trials). A t -statistic value close to 0 indicates linear scaling of cue-evoked
1807 responses with reward probability; a t -statistics value lower or greater than 0 indicates a concave
1808 or convex function of cue-evoked responses against reward probability, indicative of a
1809 pessimistic or an optimistic bias, respectively.

1810

1811



1812

1813 **Figure 5 | Model 2 cannot explain optimistic biases in behavior and cue-evoked dopamine**
 1814 **responses of habenula lesioned animals.**

1815 **a.** Possible changes in habenula lesion mice that could explain optimistic biases based on Model
 1816 2. At the level of the population dopamine responses, an optimistic bias can be caused by an
 1817 increase in the slope of the average reward responses to positive RPEs and/or a decrease in the
 1818 slope of the average reward responses to negative RPEs.

1819 **b.** At the level of the distribution of individual dopamine neuron responses, an optimistic bias
 1820 can be caused by an overall increase in the mean of the distribution of asymmetric scaling factors
 1821 (τ_i), computed from each individual neuron response function.

1822 **c.** Observed reward responses as a function of RPEs, averaged across the population of dopamine
1823 neurons for the control and lesion groups (left, mean \pm s.e.m across neurons, $n = 45$ control
1824 group, $n = 44$ lesion group). There was a significant decrease in the reward responses for the
1825 50% cue (U-statistic = 3.726, $P = 0.000195$, two-sided Mann-Whitney U-test) and 90% cue (U-
1826 statistic = 2.987, $P = 0.00281$, two-sided Mann-Whitney U-test), and for the omission responses
1827 for the 90% cue (U-statistic = -4.940, $P < 10^{-4}$, two-sided Mann-Whitney U-test). Distribution of
1828 asymmetric scaling factors (τ), computed from the average response function over the recorded
1829 neurons for the control and lesion groups (right). The distributions are the result of bootstrapping
1830 by randomly sampling neurons in 5,000 iterations. The distribution of differences between the
1831 obtained asymmetric scaling factors ($\tau_{lesion} - \tau_{control}$) was not significantly larger than zero (5th
1832 percentile = -0.1605).

1833 **d.** Distribution of asymmetric scaling factors (τ_i), computed from each individual neuron
1834 response function for the control and lesion groups. Each dot represents a single neuron ($n = 45$
1835 control group, $n = 44$ lesion group), and the neurons were sorted by asymmetric scaling factors
1836 (τ_i). The means were not significantly different (right) (t -statistic = 0.3277, $P = 0.627$, t -test).

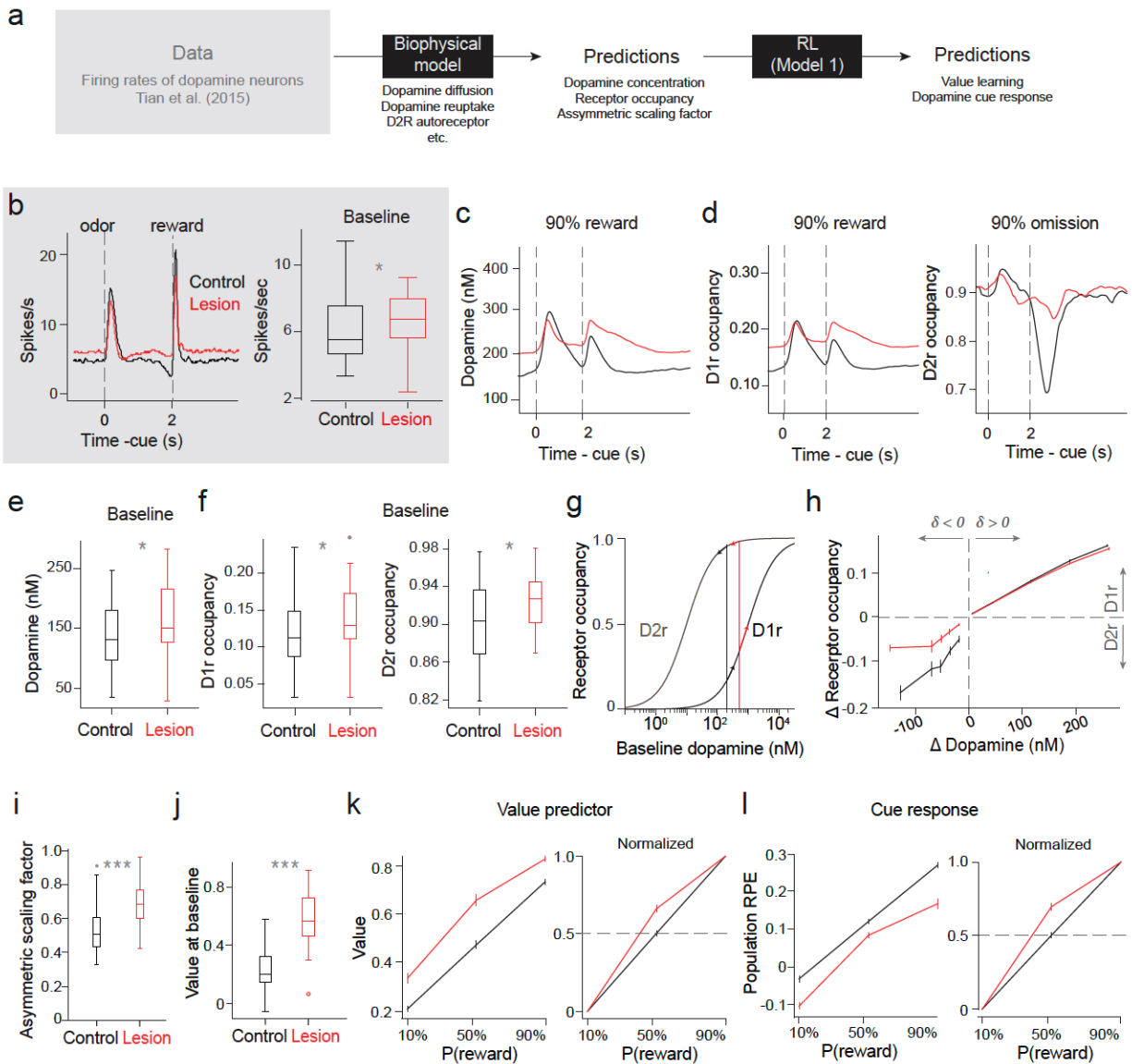
1837 **e.** Value predictions based on a TD learning model trained using the assumptions of Model 2 and
1838 the asymmetric scaling factors derived from the data. The model did not show any optimistic
1839 bias in the value predictors of the model trained with the lesion-derived asymmetric scaling
1840 factors.

1841 **f.** TD errors at cue show no signs of an optimistic bias in the model trained with the lesion-
1842 derived asymmetric scaling factors.

1843 Centre of box plot shows the median; edges are 25th and 75th percentiles; and whiskers are the
1844 most extreme data points not considered as outliers.

1845

1846



1847

1848 **Figure 6 | Model 1 can account for optimistic biases in reward-seeking behavior and cue-**
 1849 **evoked dopamine responses.**

1850 **a.** Schematic of the analysis. A biophysical model was used to predict dopamine concentrations,
 1851 receptor occupancies, and value learning based on firing rates of dopamine neurons recorded in
 1852 Tian et al. (2015).

1853 **b.** Average firing rates of dopamine neurons across the population for the control and lesion
 1854 groups (left, $n = 45$ control group, $n = 44$ lesion group). Baseline firing rates were significantly
 1855 greater in the lesion compared to the control group (right) (U-statistic = -2.429 , $P = 0.0151$,
 1856 single-sided Mann-Whitney U-test).

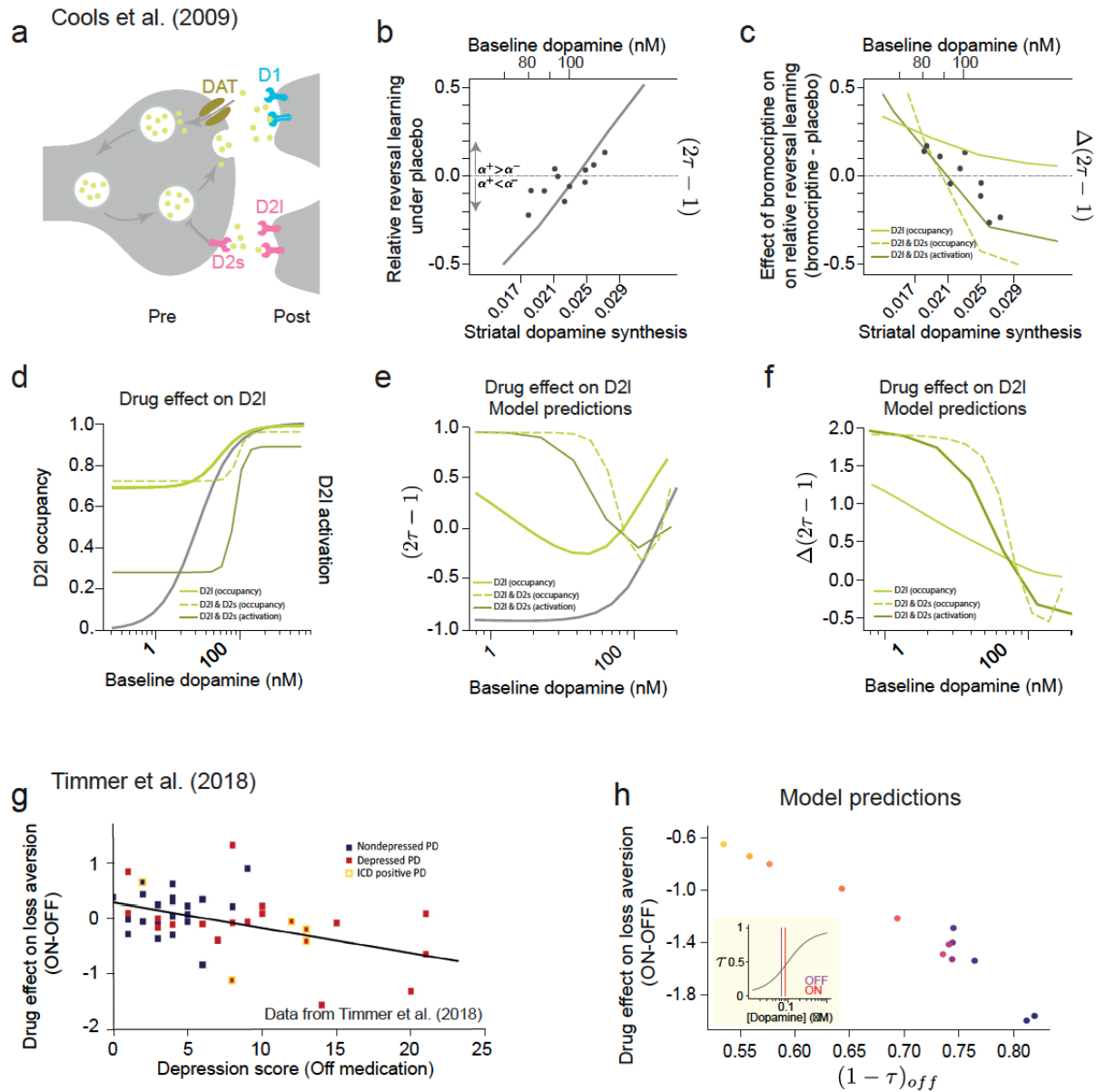
- 1857 **c.** Dopamine concentrations predicted from the firing rates of dopamine neurons based on the
1858 biophysical model of dopamine. Predictions for 90% reward trials are shown.
- 1859 **d.** Receptor occupancies predicted by the same biophysical model. Predictions for rewarded
1860 (left) and reward omission (right) trials in 90% reward trials are shown separately for D1R (left)
1861 and D2R (right), respectively ($n = 45$ control group, $n = 44$ lesion group).
- 1862 **e.** Mean dopamine concentrations at baseline predicted by the model (U-statistic = -2.109 , $P =$
1863 0.0175 , single-sided Mann-Whitney U-test).
- 1864 **f.** Mean receptor occupancies at baseline predicted by the model ($n = 45$ control group, $n = 44$
1865 lesion group). There is a significant increase in occupancies for both the D1R and D2R in the
1866 lesion compared to the control group (U-statistic = -2.1664 , $P = 0.0151$, U-statistic = -2.1328 , $P =$
1867 0.0165 for D1R and D2R respectively, single-sided Mann-Whitney U-test).
- 1868 **g.** Schematic showing the model predicted changes in dopamine concentrations and receptor
1869 occupancies for the control (black) and lesion (red) groups. The arrows depict the increase or
1870 decrease in occupancy for a positive or negative dopamine transient of a fixed magnitude.
- 1871 **h.** Changes in receptor occupancy as a function of dopamine transients predicted by the model.
1872 The slope for the positive and negative domains correspond to the receptor sensitivities of D1R
1873 and D2R (α^+ , α^-), respectively.
- 1874 **i.** Asymmetric scaling factors derived from the receptors' sensitivities for the control and lesion
1875 groups (i.e., τ in model 1, $n = 45$ control group, $n = 44$ lesion group). There was a significant
1876 increase in the lesion group with respect to controls (U-statistic = -12.205 , $P < 1.0 \times 10^{-6}$, single-
1877 sided Mann-Whitney U-test). Note that the increase in the asymmetry was driven mainly due to
1878 decreases in D2R sensitivity (panel h).
- 1879 **j.** Value predictions at baseline in a TD learning model trained with the receptor sensitivities
1880 derived from the biophysical. There was a significant increase in the value predictors at baseline
1881 in the model using the lesion group's derived parameters with respect to control. controls (t -
1882 statistic = -6.417 , $P < 1.0 \times 10^{-6}$, t -test).
- 1883 **k.** Value predictions at convergence of a TD learning model trained using the assumptions of
1884 Model 1 and the asymmetric scaling factors derived from receptors' sensitivities predicted by the
1885 biophysical model. The model led to a significant increase in the value predictions for all cues
1886 (U-statistic = -4.690 , $P < 1.0 \times 10^{-4}$, U-statistic = -4.734 , $P < 1.0 \times 10^{-4}$, U-statistic = -4.602 , $P <$

1887 1.0×10^{-4} , single-sided Mann-Whitney U-test, for the 10%, 50% and 90% reward probability
1888 cues) and an optimistic bias in the normalized value prediction to the 50% reward probability cue
1889 (t -statistic = -5.576 , $P < 1.0 \times 10^{-4}$, t -test) in accordance with the anticipatory licking observed in
1890 the data.

1891 **I. Predicted cue responses.** There is an overall decrease in RPEs in lesioned animals (left) due to
1892 an increase in the baseline (pre-cue) value prediction (U-statistic = 4.932 , $P < 1.0 \times 10^{-5}$, U-
1893 statistic = -3.658 , $P = 0.00025$, U-statistic = 4.734 , $P < 1.0 \times 10^{-4}$, single-sided Mann-Whitney
1894 U-test for the 10%, 50% and 90% reward probability cues), which is consistent with the
1895 decreases in the absolute magnitudes of dopamine cue-evoked responses in the lesion group (Fig.
1896 4c). The normalized TD errors at for the 50% reward probability cue show signs of an optimistic
1897 bias (U-statistic = -4.624 , $P < 1.0 \times 10^{-4}$, single-sided Mann–Whitney U-test).

1898 Centre of box plot shows the median; edges are 25th and 75th percentiles; and whiskers are the
1899 most extreme data points not considered as outliers.

1900



1901

1902 **Figure 7 | Model 1 predicts asymmetric learning rates in healthy humans given inter-**
 1903 **individual differences in baseline dopamine, and in Parkinson’s disease patients given**
 1904 **inter-individual differences in depressive-like symptoms.**

1905 **a.** Schematic of the events occurring at dopaminergic axon terminal. Pre- and post-synaptic sites
 1906 predominantly express D2s (short) and D2l (long) subtypes, respectively.

1907 **b.** “Relative reversal learning (RRL)” under placebo conditions as a function of dopamine striatal
 1908 synthesis capacity measured with PET radio imaging (black dots, left y-axis, bottom x-axis).

1909 Figure taken from Cools et al. (2009)²⁴. Positive values of RRL indicate a bias favoring learning
1910 from gains relative to losses, and vice versa for negative values of RRL. There was a positive
1911 relationship between RRL and dopamine synthesis capacity. Model 1 predictions of RRL ($2\tau -$
1912 1 in Model 1) as a function of baseline dopamine using the receptors occupancy curve,
1913 recapitulate the positive relationship shown in the results from Cools et al. (2009)²⁴ (gray line,
1914 right y-axis, top x-axis).

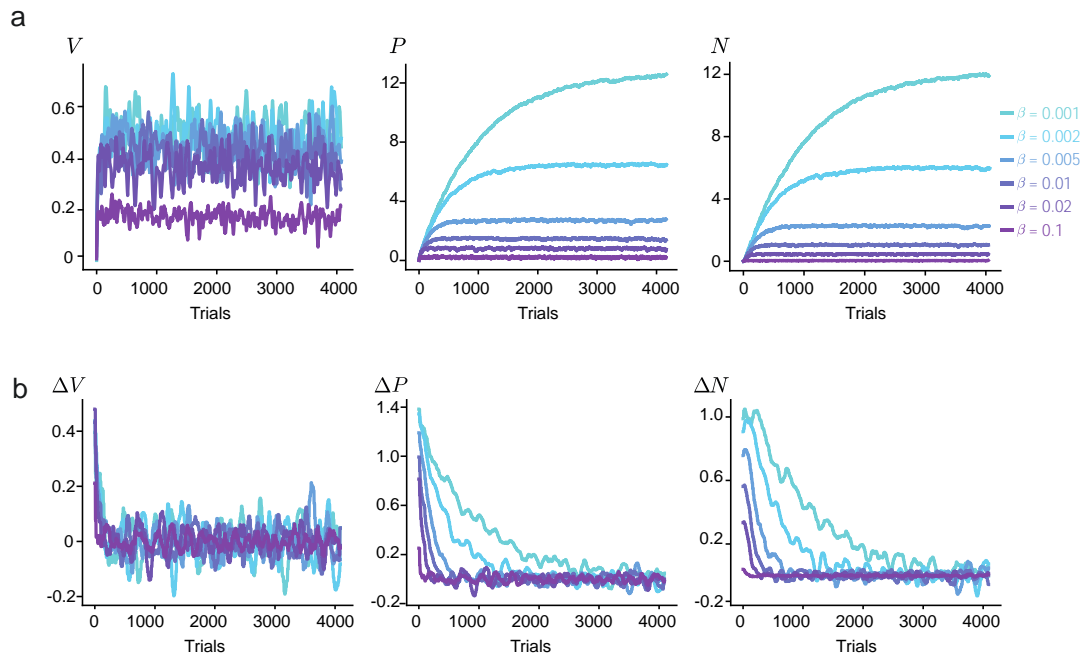
1915 **c.** The change in RRL induced by bromocriptine was negatively correlated with striatal
1916 dopamine synthesis capacity. Figure from Cools et al. (2009)²⁴ (black dots, left y-axis, bottom x-
1917 axis). Model 1 recapitulates qualitatively the effect of bromocriptine in RRL, equivalent to
1918 $\Delta(2\tau - 1)$. The solid light green line represents the $\Delta(2\tau - 1)$ when considering bromocriptine's
1919 effect on *D2l occupancy only*; the dashed line represents the $\Delta(2\tau - 1)$ when both *D2l and D2s*
1920 *occupancy* was considered; and the dark green line represents the $\Delta(2\tau - 1)$ when both *D2l and*
1921 *D2s activation* was considered (this includes the fact that bromocriptine is a partial agonist for
1922 the D2l and D2s receptors). The curves were obtained by imposing a concentration of $10^{0.8}$ nM
1923 of bromocriptine in the biophysical model.

1924 **d.** Receptor occupancy curves for the D2l receptors at baseline (grey line) and in the presence of
1925 $10^{0.8}$ nM of bromocriptine: Solid light green line corresponds to considering bromocriptine's
1926 effects on D2l receptors occupancy alone; dashed line, corresponds to considering
1927 bromocriptine's effects on both D2l and D2s receptors; solid dark green line corresponds to
1928 considering bromocriptine's effect on the activation curves of both D2l and D2s receptors. The
1929 binding of the drug to the D2l receptors alone causes an increase in occupancy. This happens to a
1930 larger extent when starting from a low dopamine level at baseline than in high dopamine levels.
1931 The binding of the drug to D2s receptors in addition to D2l receptors causes a rightward shift in
1932 the curves. The activation levels are lower than 1 even at the drug levels where occupancy is
1933 close to 1, due to the lower efficiency of bromocriptine in receptor activation (Methos 4.1). See
1934 Extended Figure 6 and 7 for the effect of changing bromocriptine's concentration and efficiency
1935 of activation.

1936 **e.** Same as in panel **d.** but now reporting $(2\tau - 1)$ calculated from the D2l receptor's occupancy
1937 and activation curves. An increase in $2\tau - 1$ happens to a larger extent when starting from a low
1938 dopamine level than from high dopamine level.

1939 **f.** Same as in panel **d.** but now reporting $\Delta(2\tau - 1)$ calculated from the D2l receptor's occupancy
1940 and activation curves. Model 1 recapitulates qualitatively the effect of bromocriptine on RRL.
1941 **g.** The effect of PD medication (L-DOPA) on loss aversion is negatively correlated with their
1942 off-medication depression score. Figure from Timmer et al. (2018)⁷¹.
1943 **h.** Model 1 recapitulates qualitatively the effect of PD medication in loss aversion. We assumed
1944 that the asymmetry in favor of learning from losses relative to gains $(1 - \tau)_{off}$ scales with the
1945 baseline dopamine levels. Given this, we derived a distribution of off-medication baseline
1946 dopamine levels centered around the mean $(1 - \tau)_{off}$ derived from the data of Timmer et al.
1947 (2018)⁷¹ (see methods 0). We then imposed a fixed increase in baseline dopamine to simulate L-
1948 DOPA effects. We derived the new loss-aversion parameter $(1 - \tau)_{on}$ at the shifted baseline
1949 dopamine levels. The y-axis shows the change in loss aversion for each sample of the
1950 distribution of baseline dopamine levels. If the off-medication depression score is correlated with
1951 $(1 - \tau)_{off}$ then model would predict the result in Timmer et al. (2018)⁷¹.
1952 PD: Parkinson's disease.
1953

1954 **Extended data Figures**



1955

1956 **Extended Data Fig. 1 | Variables of model 1 show convergence irrespective of the value of**
1957 **the decay factor.**

1958 **a.** Value predictors V (left), P population (middle) and N population (right) across trials of
1959 training for an RL agent of model 1. Color of lines denotes the value of the decay factor (β) in
1960 the update rules for the P and N populations. Colormap is the same for all figures (left). All the
1961 model variables show convergence for every value of the decay factor β .

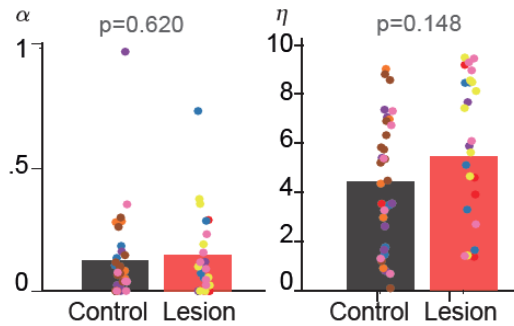
1962 **b.** Difference in the variables estimates between consecutive trials of training, for the value
1963 predictors V (left, ΔV), P population (middle, ΔP) and N population (right, ΔN). All the variables
1964 show convergence for every value of the decay factor β (shown as a ΔV , ΔP , ΔN equal to zero).

a TD learning

$$\delta = r - V$$

$$V \leftarrow V + \alpha \cdot \delta$$

$$\text{Licking} = \eta \cdot V$$

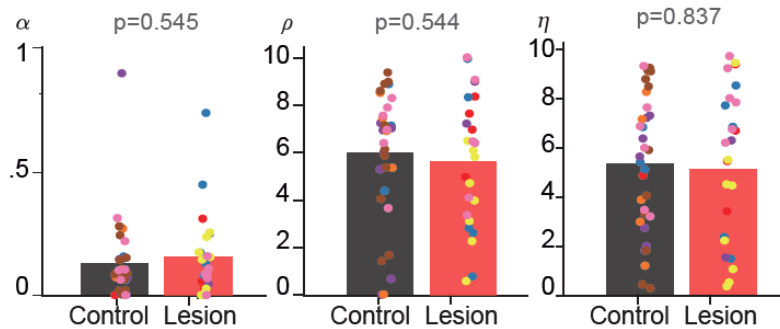


b TD learning with reward sensitivity

$$\delta = \rho \cdot r - V$$

$$V \leftarrow V + \alpha \cdot \delta$$

$$\text{Licking} = \eta \cdot V$$



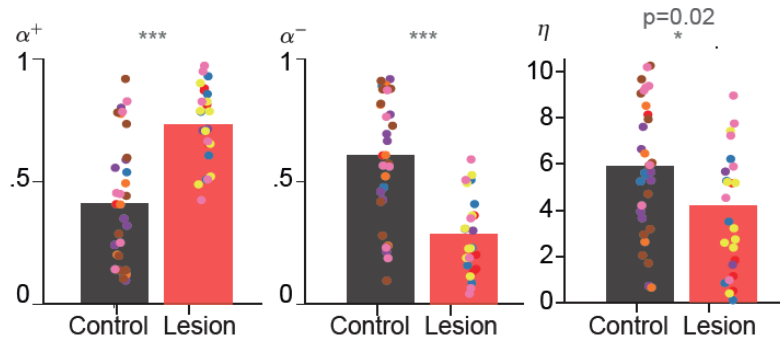
c Risk sensitive reinforcement learning

$$\delta = r - V$$

$$V \leftarrow V + \alpha^+ \cdot \delta \quad \text{if } \delta > 0$$

$$V \leftarrow V + \alpha^- \cdot \delta \quad \text{if } \delta < 0$$

$$\text{Licking} = \eta \cdot V$$



1965

1966

1967 **Extended Data Fig. 2 | RL model fits to the trial-by-trial anticipatory licking responses.**

1968 **a.** TD learning fits reveal no significant difference across groups in the learning rate (left, U-
 1969 statistic = -4.954 , $P = 0.620$, two-sided Mann-Whitney U-test) nor in the linear mapping between
 1970 value predictions and anticipatory licking responses (U-statistic = -1.445 , $P = 0.148$, two-sided
 1971 Mann-Whitney U-test).

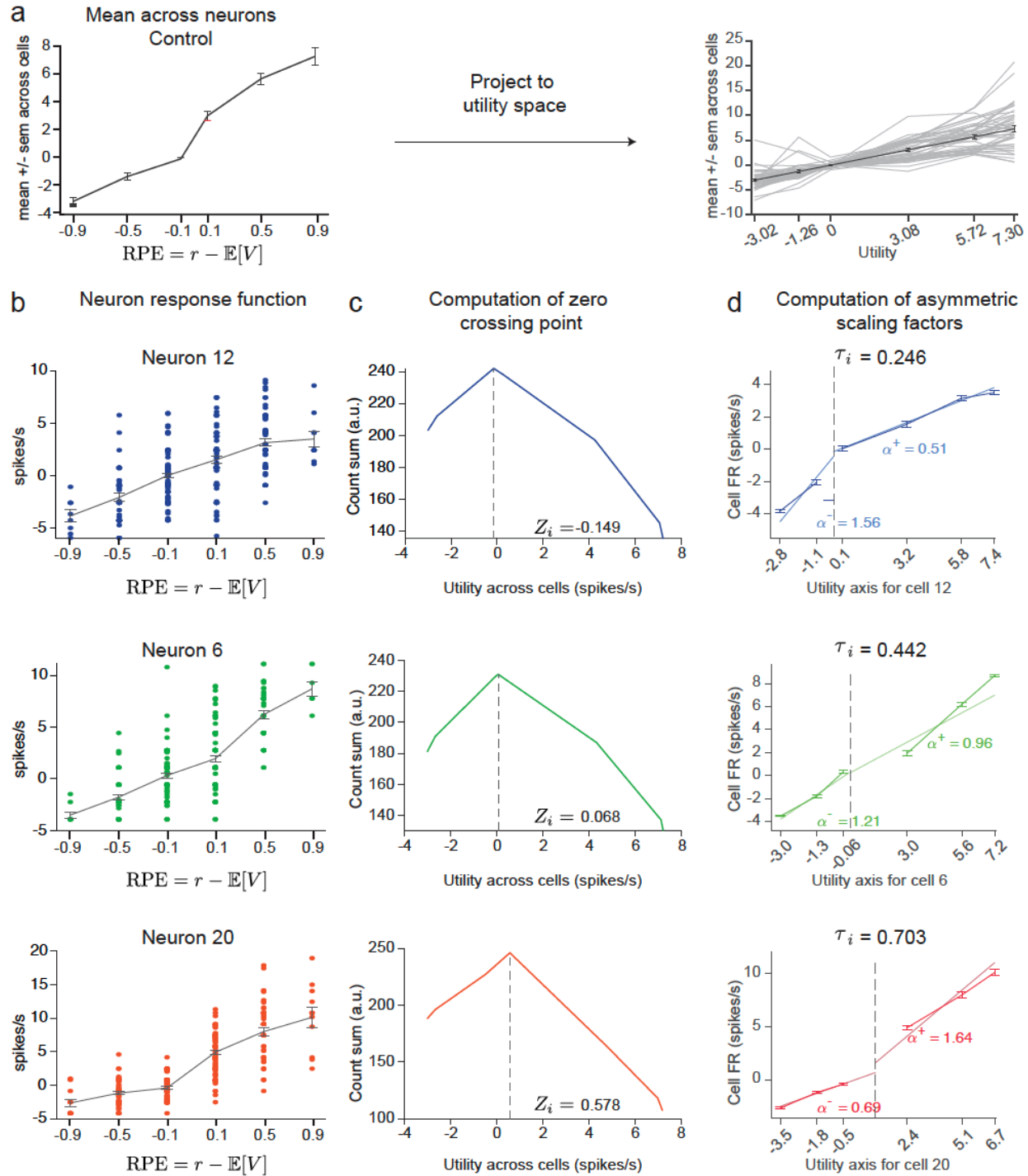
1972 **b.** Model fits of TD learning with reward sensitivity reveal no difference across groups in the
 1973 learning rate (left, U-statistic = 0.206 , $P = 0.836$, two-sided Mann-Whitney U-test) nor in the

1974 linear mapping between value predictions and anticipatory licking responses (middle, U-statistic
1975 = -0.7844 , $P = 0.4327$, two-sided Mann-Whitney U-test), nor in the reward sensitivity (right) (U-
1976 statistic 0.545 , $P = 0.605$, two-sided Mann-Whitney U-test).

1977 c. Model fits of TD learning with asymmetric learning rates for positive vs negative RPEs. This
1978 model reveals a significant difference across groups in the asymmetry between α^+ and α^- (U-
1979 statistic = -4.678 , $P < 1.0 \times 10^{-5}$, two-sided Mann-Whitney U-test) and a small but significant
1980 difference between the linear mapping between value predictions and anticipatory licking
1981 responses (right, U-statistic = 2.33 , $P = 0.02$, two-sided Mann-Whitney U-test).

1982

1983



1984

1985 **Extended Data Fig. 3 | Signatures of distributional reinforcement learning model are**
 1986 **preserved after habenula lesions.**

1987 **a.** RPE -evoked responses at outcome as a function of the theoretical RPE for each trial type. The
 1988 figure shows the average response function across neurons from the control group. The

1989 computation of zero-crossing points and asymmetric scaling factors is carried out in the ‘utility
1990 space’(the average response function) as in ³⁶ to account for response nonlinearities.

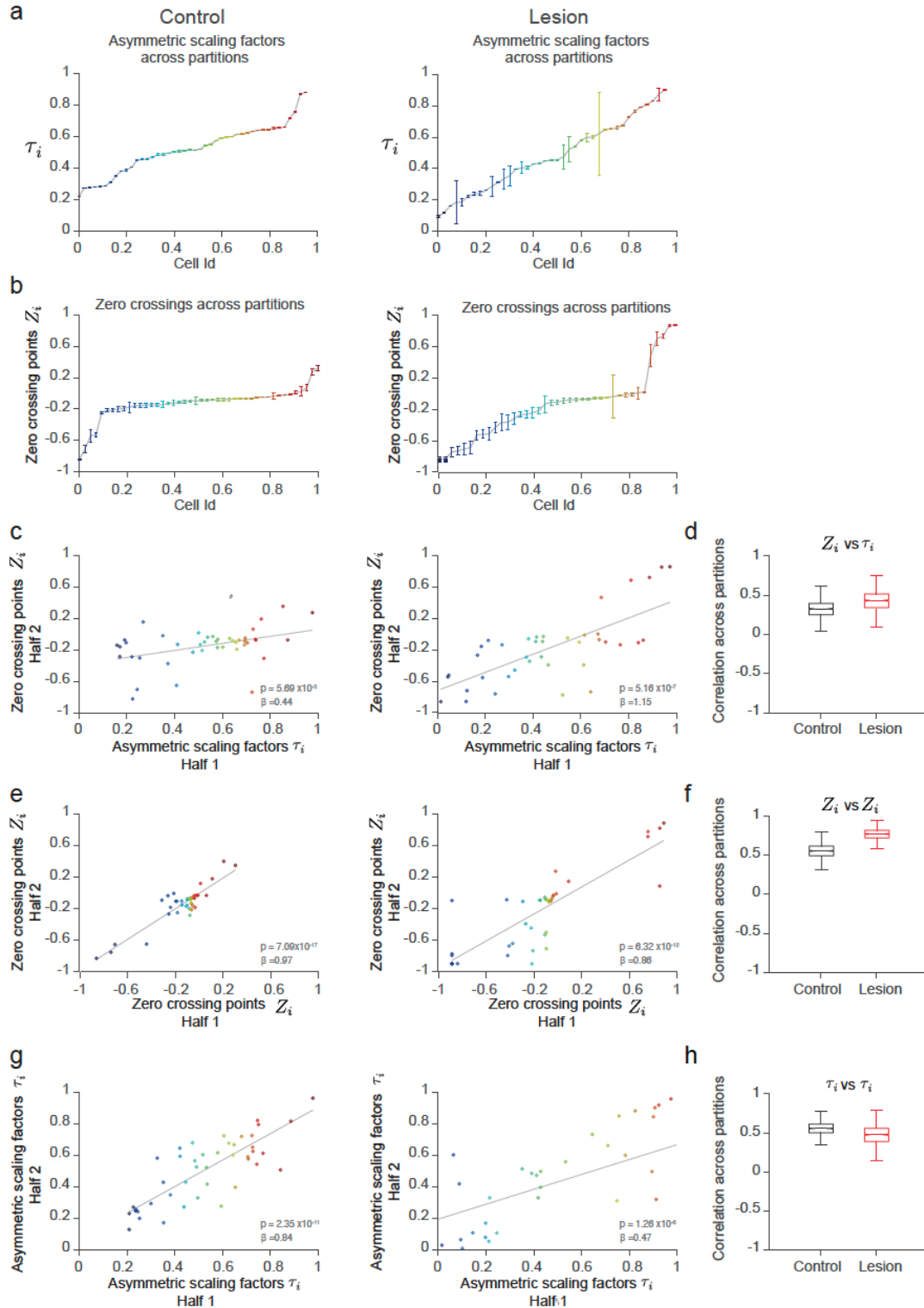
1991 **b.** Example of the response function of 3 dopamine neurons from the control group ordered by
1992 their asymmetric scaling factors: pessimistic, neutral and optimistic, from top to bottom.

1993 **c.** Computation of zero-crossing points for the neurons in B. The reversal points for each cell (Z_i)
1994 were defined as the point in utility space that maximized the number of positive responses to
1995 points greater than Z_i plus the number of negative responses to points less than Z_i . The y-axis
1996 shows the sum of responses below and above each point in the utility space. The zero-crossing
1997 point is shown as the maxima in this function.

1998 **d.** Computation of asymmetric scaling factors for the neurons in c. Here, the responses functions
1999 in b have been projected to the utility space in A and realigned according to their zero-crossing
2000 points. The asymmetric learning rates (α^+ , α^-) are taken to be the slopes of these response
2001 functions.

2002

2003



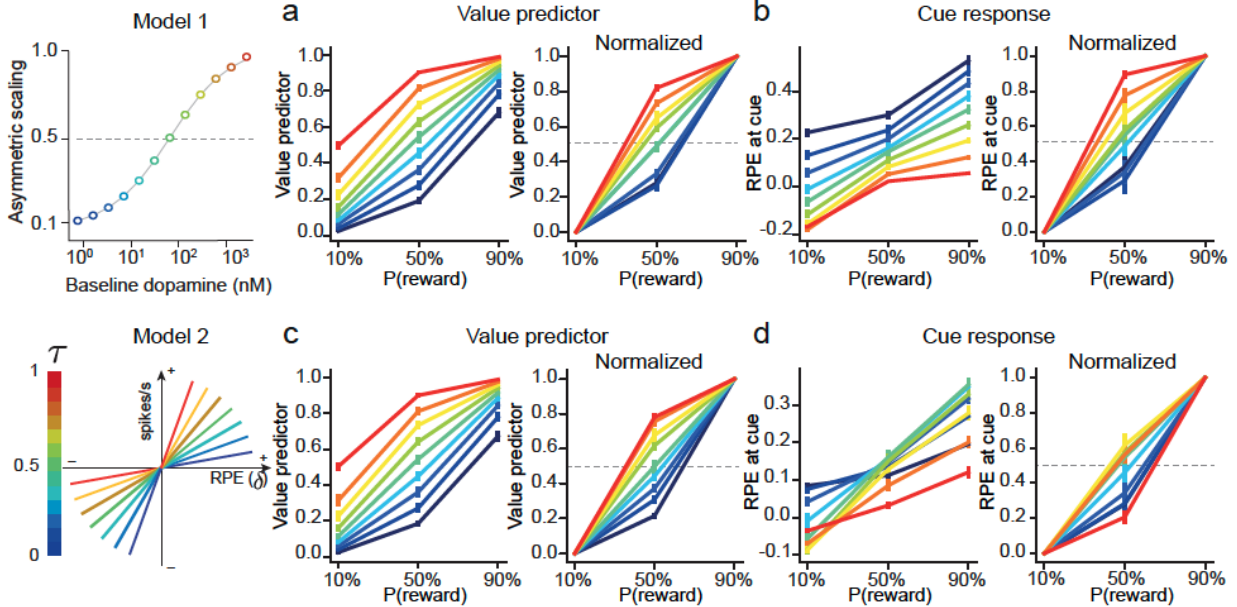
2004

2005 **Extended Data Fig. 4 | Distributional reinforcement learning variables from the Habenula**

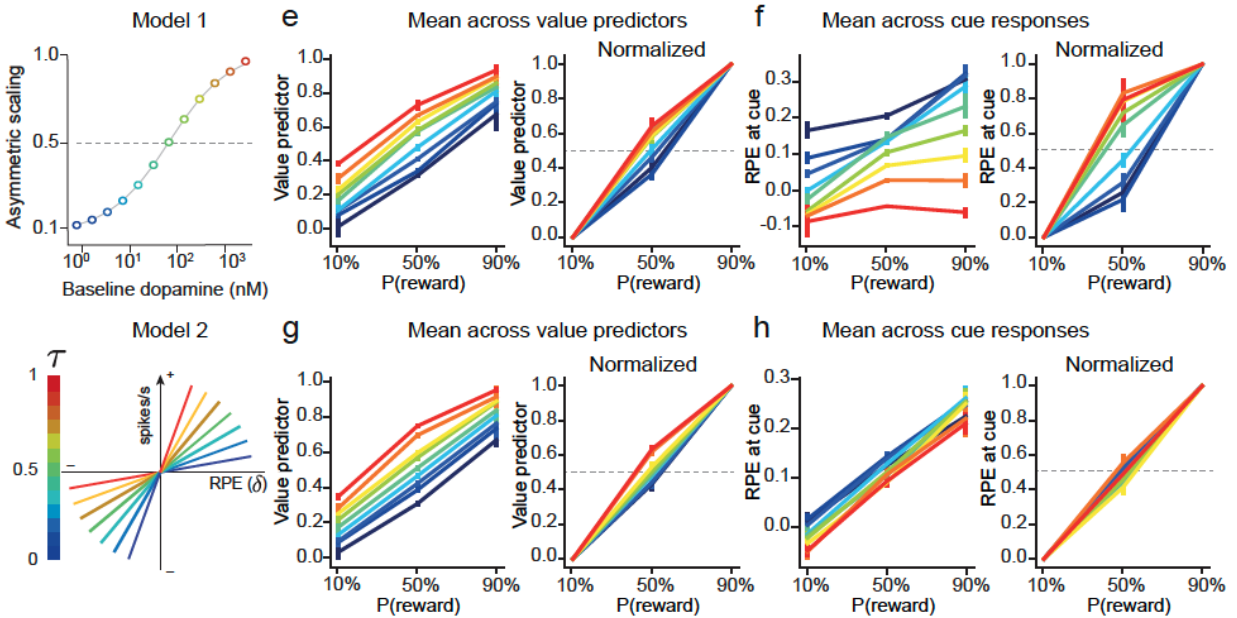
2006 **lesion dataset.**

- 2007 **a.** Distribution of asymmetric scaling factors for the dopamine neurons from the control (left)
2008 and lesion (right) groups. The error bars were derived by randomly sampling trials to compute
2009 the asymmetric scaling factors for 1000 iterations.
- 2010 **b.** Distribution of zero crossing points for the control (left) and lesion(right) groups. The error
2011 bars were derived as in a.
- 2012 **c.** Correlation of asymmetric scaling factors (x-axis) and zero-crossing points (y-axis) computed
2013 on disjoint halves of trials for an example partition.
- 2014 **d.** Distribution of correlation coefficients between asymmetric scaling factors (x-axis) and zero-
2015 crossing points (y-axis) across disjoint halves of trials for 1000 partitions for the control and
2016 lesion groups.
- 2017 **e.** Correlation between zero-crossing points computed on disjoint halves of trials for an example
2018 partition.
- 2019 **f.** Distribution of correlation coefficients between zero-crossing points computed on disjoint
2020 halves of trials for 1000 partitions for the control and lesion groups.
- 2021 **g.** Correlation between asymmetric scaling factors computed on disjoint halves of trials for an
2022 example partition.
- 2023 **h.** Distribution of correlation coefficients between asymmetric scaling factors computed on
2024 disjoint halves of trials for 1000 partitions for the control and lesion groups.
- 2025
- 2026

TD learning



Distributional TD learning

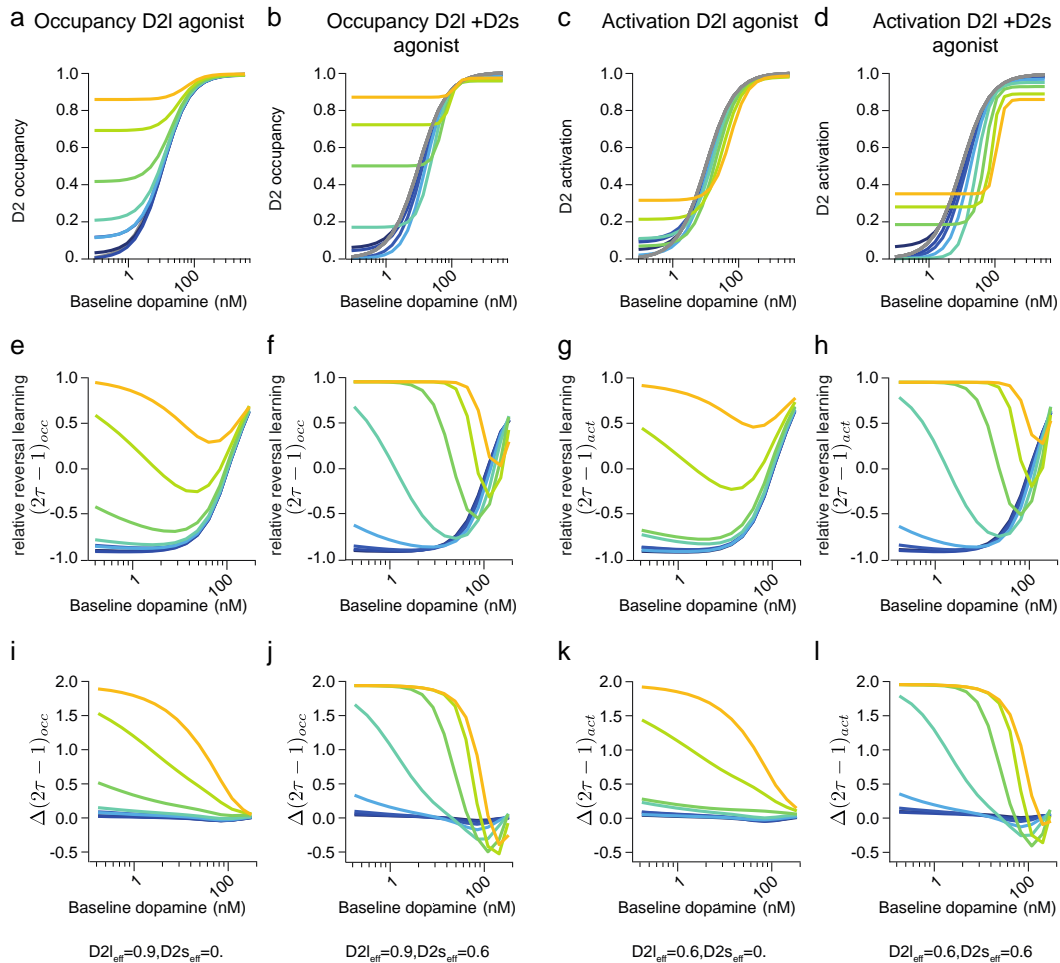


2027

2028 **Extended Data Fig. 5 | Biases in cue-evoked responses in the Habenula lesion data cannot**
 2029 **be explained by asymmetric scaling of RPEs (Model 2).**

2030 **a.** Value predictors derived from model 1 with TD learning for a set of baseline dopamine levels
 2031 (colormap). The optimistic and pessimistic biases are present.

- 2032 **b.** Cue responses derived from model 1 with TD learning for a set of baseline dopamine levels
2033 (colormap). The optimistic and pessimistic biases are revealed when the responses are
2034 normalized.
- 2035 **c.** Value predictors derived from model 2 with TD learning for a set of baseline dopamine levels
2036 (colormap). The optimistic and pessimistic biases are present.
- 2037 **d.** Cue responses derived from model 2 with TD learning for a set of baseline dopamine levels
2038 (colormap). The optimistic and pessimistic biases are absent in both the normalized and the raw
2039 TD errors.
- 2040 **e.** Mean across the distribution of value predictors derived from model 1 with distributional TD
2041 learning for a set of baseline dopamine levels (colormap). The optimistic and pessimistic biases
2042 are present.
- 2043 **f.** Mean across the distribution of cue responses derived from model 1 with distributional TD
2044 learning for a set of baseline dopamine levels (colormap). The optimistic and pessimistic biases
2045 are revealed when the responses are normalized.
- 2046 **g.** Mean across the distribution of value predictors derived from model 2 with distributional TD
2047 learning for a set of baseline dopamine levels (colormap). The optimistic and pessimistic biases
2048 are present.
- 2049 **h.** Mean across the distribution of cue responses derived from model 2 with distributional TD
2050 learning for a set of baseline dopamine levels (colormap). The optimistic and pessimistic biases
2051 are absent in both the normalized and the raw TD errors.
- 2052
- 2053



2054

2055 **Extended Data Fig. 6 | Model 1 predicts asymmetric learning rates and the effect of**
 2056 **bromocriptine in healthy humans given inter-individual differences in baseline dopamine.**

2057 **a.** Occupancy curves for the D2I receptors at baseline (grey line) and when considering
 2058 bromocriptine's effects in D2I receptors alone. The binding of the drug to the D2I receptors alone
 2059 causes an increase in the occupancy. This happens to a larger extent when starting from a low
 2060 dopamine level at baseline than in high dopamine levels.

2061 **b.** Occupancy curves for the D2I receptors at baseline (grey line) and when considering
 2062 bromocriptine's effects in both D2I and D2s receptors. The binding of the drug to D2s receptors
 2063 causes a rightwards shifts in the curves.

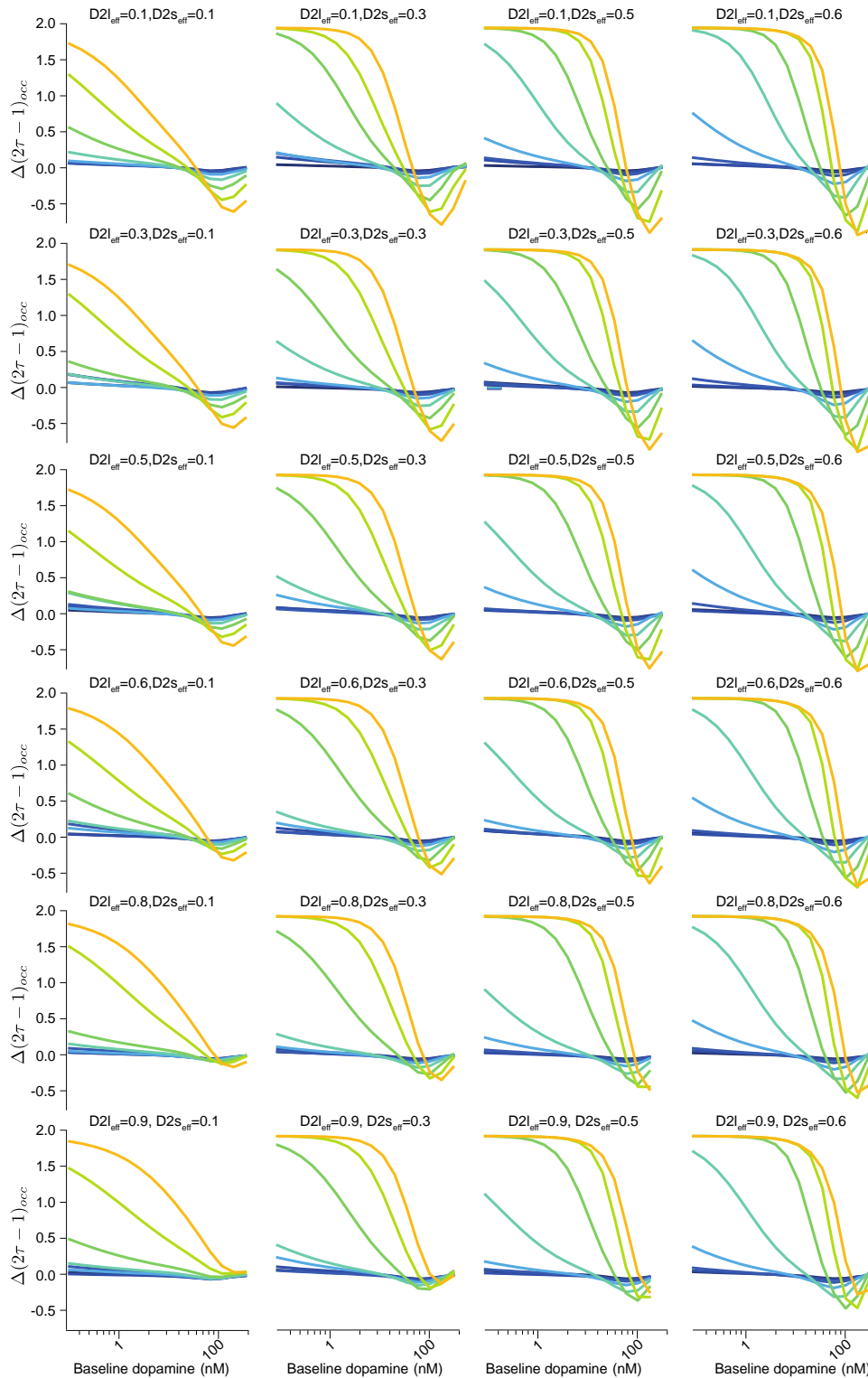
- 2064 **c.** Activation curves for the D2l receptors at baseline (grey line) and when considering
2065 bromocriptine's effects in D2l receptors alone, including the *partial* quality of the agonism of
2066 this drug on the receptor (i.e., efficiency <1).
- 2067 **d.** Activation curves for the D2l receptors at baseline (grey line) and when considering
2068 bromocriptine's effects in both D2l and D2s receptors, including the *partial* quality of the
2069 agonism.
- 2070 **e.** Relative reversal learning (RRL) calculated as $2\tau - 1$ in model 1, as a function of baseline
2071 dopamine (x-axis) and drug concentration (color) using the receptors occupancy curve,
2072 considering only bromocriptine's effect in D2l receptors.
- 2073 **f.** Relative reversal learning (RRL) calculated as $2\tau - 1$ in model 1, as a function of baseline
2074 dopamine (x-axis) and drug concentration (color) using the receptors occupancy curve,
2075 considering bromocriptine's effect in both D2l and D2s receptors.
- 2076 **g.** $2\tau - 1$ in model 1, as a function of baseline dopamine (x-axis) and drug concentration (color)
2077 using the receptors activation curve, considering only bromocriptine's effect in D2l receptors.
- 2078 **h.** $2\tau - 1$ in model 1, as a function of baseline dopamine (x-axis) and drug concentration (color)
2079 using the receptors activation curve, considering bromocriptine's effect in both D2l and D2s
2080 receptors.
- 2081 **i.** The change in $2\tau - 1$ induced by the drug at different concentrations (color) with respect to
2082 the baseline condition, as a function of baseline dopamine (x-axis). The curves represent the
2083 change when calculating $2\tau - 1$ from the occupancy curves considering only D2l binding.
- 2084 **j.** The change in $2\tau - 1$ induced by the drug at different concentrations (color) with respect to
2085 the baseline condition, as a function of baseline dopamine (x-axis). The curves represent the
2086 change when calculating $2\tau - 1$ from the occupancy curves considering both D2l and D2s
2087 binding
- 2088 **k.** The change in $2\tau - 1$ induced by the drug at different concentrations (color) with respect to
2089 the baseline condition, as a function of baseline dopamine (x-axis). The curves represent the
2090 change when calculating $2\tau - 1$ from the activation curves considering only D2l activation.
- 2091 **l.** The change in $2\tau - 1$ induced by the drug at different concentrations (color) with respect to
2092 the baseline condition, as a function of baseline dopamine (x-axis). The curves represent the

2093 change when calculating $2\tau - 1$ from the activation curves considering both D2l and D2s
2094 activation.

2095 The parameters of efficiency of activation of D2 receptors by the drug ($D2l_{eff}, D2s_{eff}$) used in
2096 each column are reported at the bottom of the figure.

2097

2098



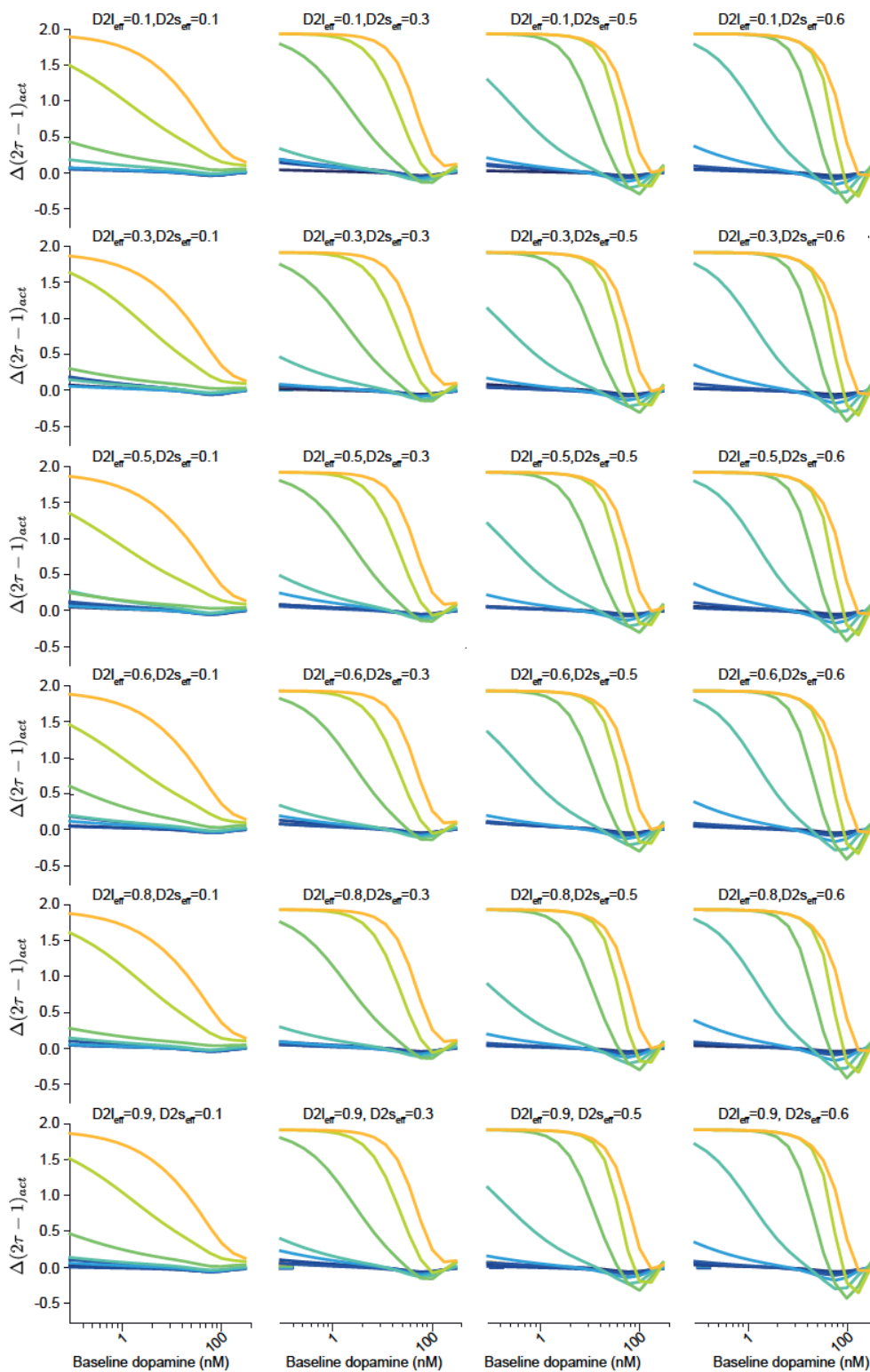
2099

2100 **Extended Data Fig. 7 | Robustness of the effect of bromocriptine in the relative reversal**
2101 **learning calculated from the occupancy curves to the choice of the drug efficiency**
2102 **parameter.**

2103 The qualitative effects on bromocriptine in the change in relative reversal learning $\Delta(2\tau - 1)$
2104 calculated from the D2 *occupancy* curves. Results hold regardless of the choice of the efficiency
2105 of the drug on D2l (rows) or D2s (columns) efficiency.

2106

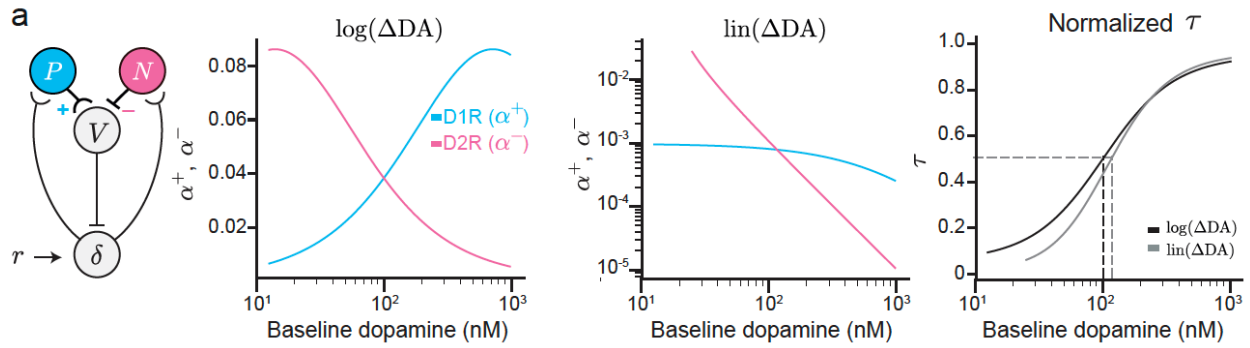
2107



2108

2109 **Extended Data Fig. 8 | Robustness of the effect of bromocriptine in the relative reversal**
2110 **learning calculated from the activation curves to the choice of the drug efficiency**
2111 **parameter.**

2112 The qualitative effects on bromocriptine in the change in relative reversal learning $\Delta(2\tau - 1)$
2113 calculated from the D2 *activation* curves. Results hold regardless of the choice of the efficiency
2114 of the drug on D2l (rows) or D2s (columns) efficiency.
2115



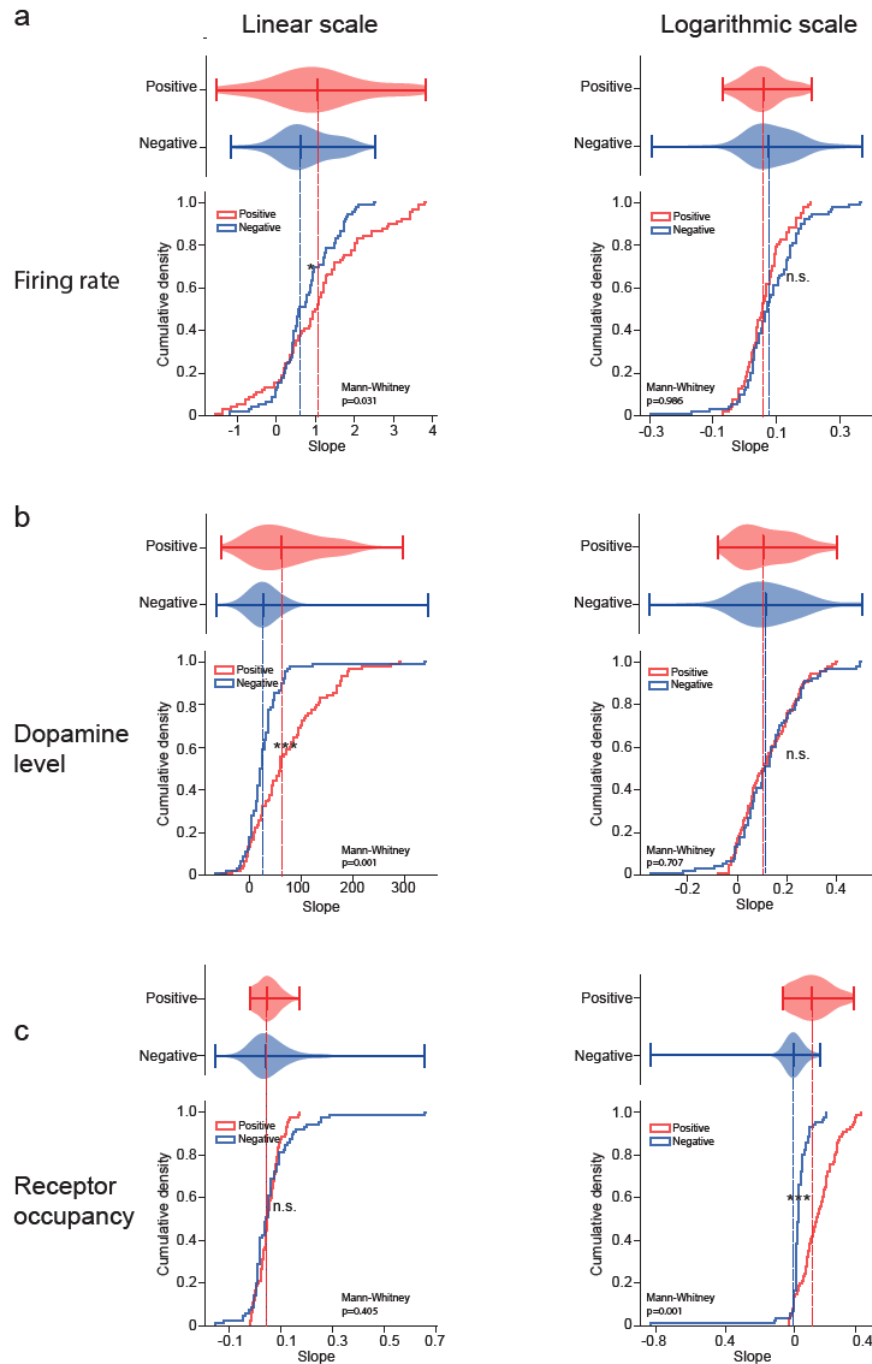
2116

2117

2118 **Extended Data Fig. 9 | The qualitative aspects of Model 1 are preserved irrespective of the**
2119 **assumption made about the changes in baseline dopamine caused by dopamine transients.**

2120 **a.** Computation of receptor sensitivities (i.e., slope of dose-occupancy curves, α^+ , α^-) assuming
2121 logarithmic (left) or linear (middle) changes in baseline dopamine induced by dopamine
2122 transients ($\log\Delta DA$, $\text{lin}\Delta DA$, respectively). The absolute magnitude of the slopes differs
2123 depending on the assumption made about the changes in baseline dopamine (logarithmic vs
2124 linear) but the asymmetric scaling factor presents only a small shift in the curve as a function of
2125 baseline dopamine (right). The qualitative aspects of the model (i.e., non-monotonic relationship
2126 of the asymmetric scaling factor with baseline dopamine) is preserved regardless on this
2127 assumption.

2128



2129

2130 **Extended Data Fig. 10 | Change in dopamine firing rates, dopamine concentration and**
2131 **receptor occupancy as a function of RPEs in the linear scale or logarithmic scale.**

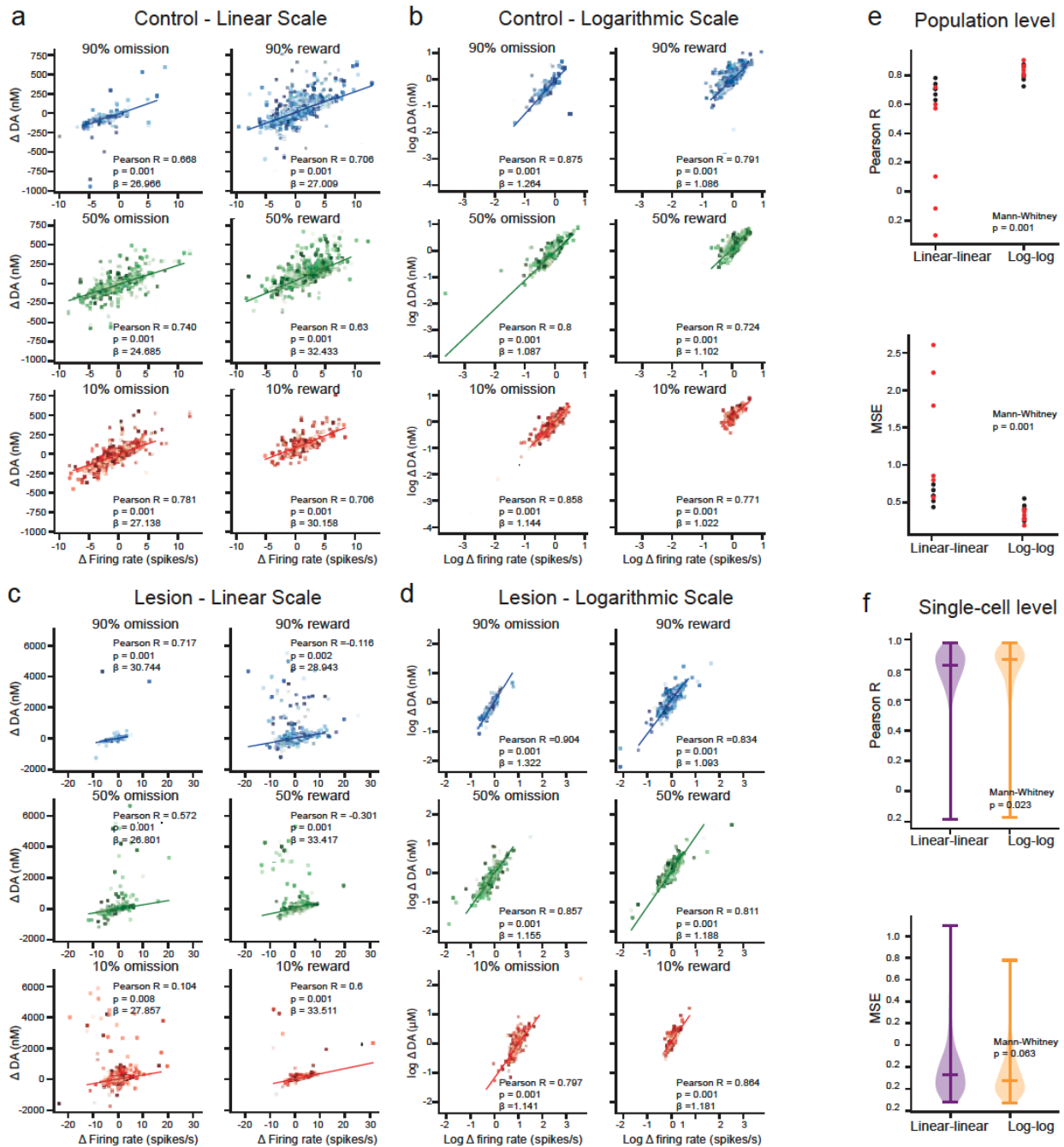
2132 **a.** Distributions of the slopes of the change in firing rate derived as a function of RPEs in the
2133 positive and negative domains, computed in the linear (left) or logarithmic (right) scale,

2134 calculated at a single neuron level. The slopes are asymmetric if considered in the linear scale,
2135 with the negative transients presenting a shallower slope than the positive ones. The slopes are
2136 symmetric if considered in the logarithmic scale.

2137 **b.** Distributions of the slopes of the change in dopamine levels derived from the biophysical
2138 model as a function of RPEs in the positive and negative domains, computed in the linear (left)
2139 or logarithmic (right) scale, calculated at a single neuron level. The slopes are again asymmetric
2140 if considered in the linear scale but symmetric if considered in the logarithmic scale.

2141 **c.** Slope of the change in receptor occupancy derived from the biophysical model for a given
2142 RPE in the positive and negative domains, computed in the linear (left) or logarithmic (right)
2143 scale, calculated at a single neuron level. The slopes are symmetric if considered in the linear
2144 scale but asymmetric if considered in the logarithmic scale.

2145



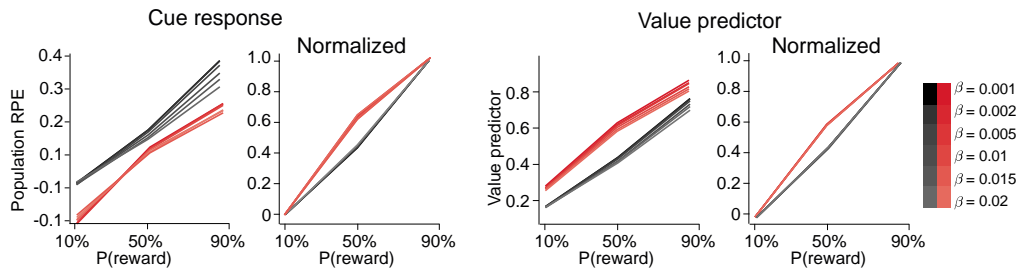
2146

2147 **Extended Data Fig. 11 | Relationship between changes in firing rates and changes in**
 2148 **dopamine concentration derived from the biophysical model.**

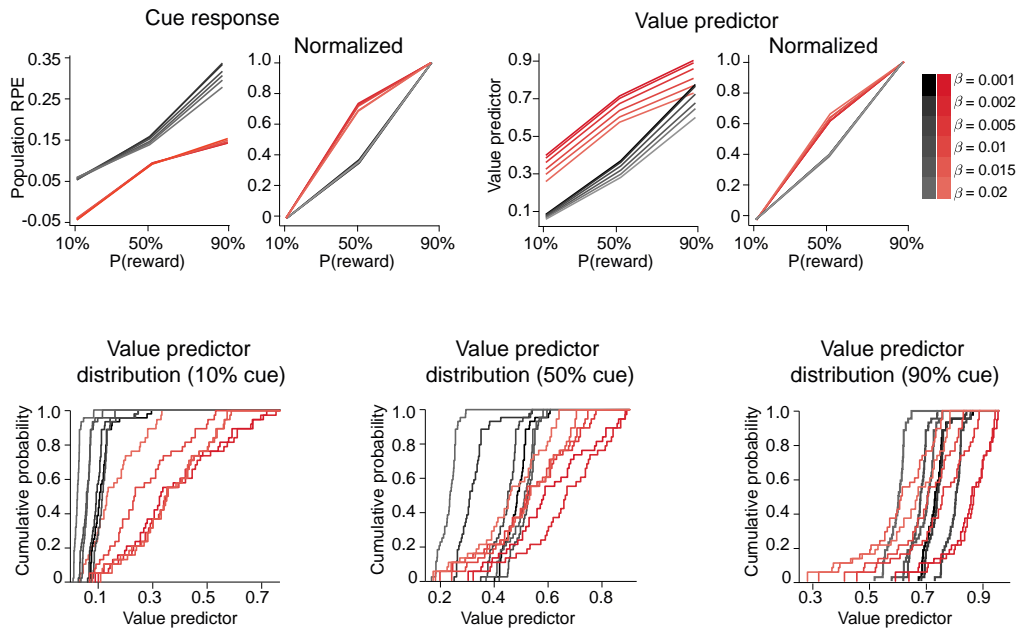
2149 **a.** Linear fits to the relationship between changes in firing rates and changes in dopamine
 2150 concentration evoked by the TD error at outcome in the linear scale for the control group. The
 2151 fits are done separately for each trial type.

- 2152 **b.** Same as a, but fits are done in the logarithmic scale.
- 2153 **c.** Same as a, but fits are done for the lesion group.
- 2154 **d.** Same as c, but fits are done in the logarithmic scale.
- 2155 **e.** Distribution of the Pearson correlation coefficients (top) and means squared error (MSE,
2156 bottom) between the predicted change in dopamine concentration by the linear regression and the
2157 ground truth derived from the biophysical model. The coefficients are derived from the fits in
2158 figures a-e, done by pooling all trials for each trial type (each point each trial type, with black for
2159 control and red for lesion group). There was a significant increase in the Pearson correlation
2160 coefficient and a near significant decrease in the MSE if the changes are considered to happen in
2161 the logarithmic scale.
- 2162 **f.** Same as E, but the linear regression fits are done for each neuron separately by pooling all
2163 trials. There was a significant increase in the single-cell distribution of Pearson correlation
2164 coefficients (top) and a significant decrease in the MSE distribution (bottom) if the changes are
2165 considered to happen in the logarithmic scale.
- 2166
- 2167

a Model 1 - No distributional TD



b Model 1 - Distributional TD



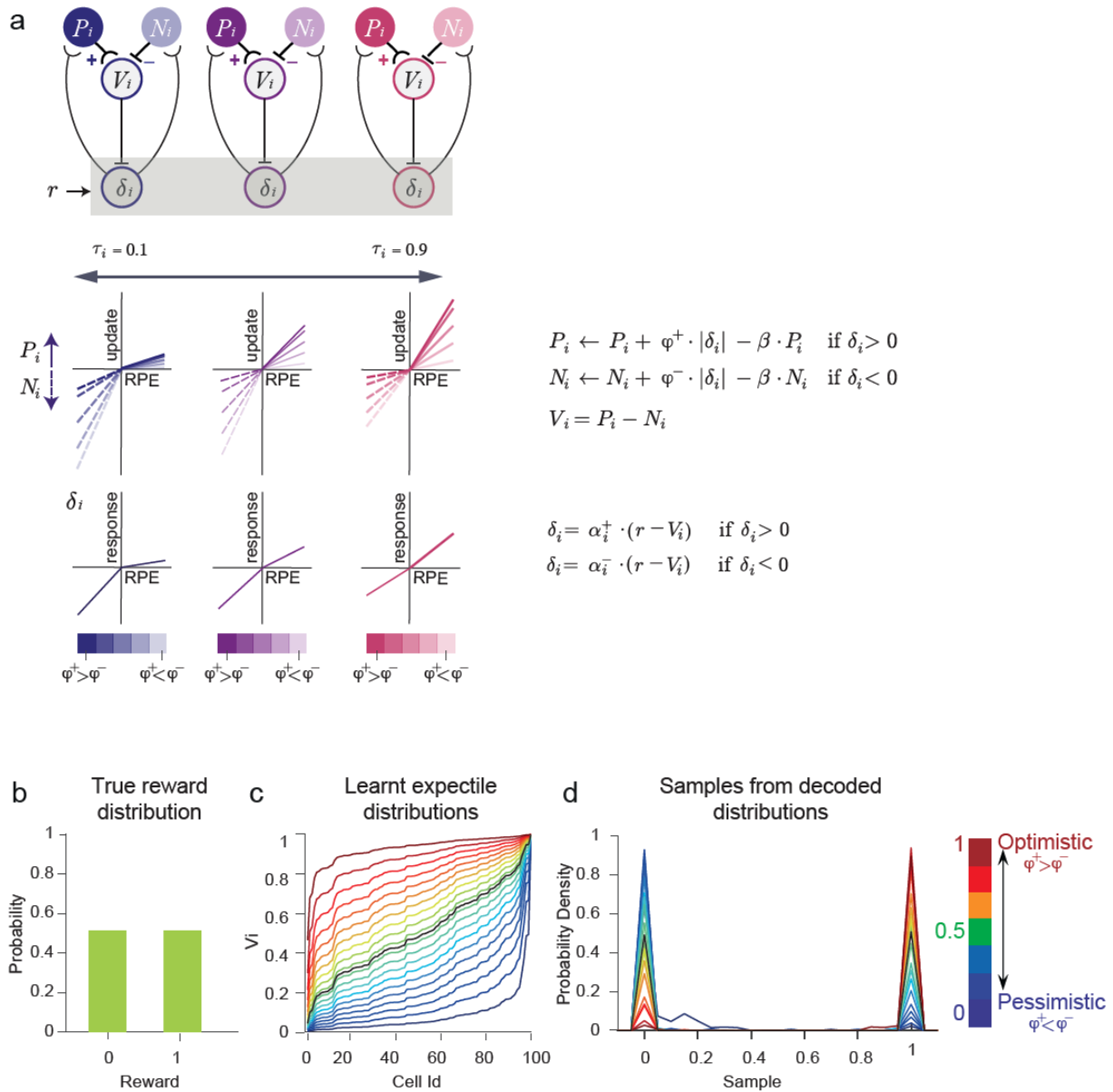
2168

2169 **Extended Data Fig. 12 | Model 1 captures signatures of the data irrespective of the choice of**
 2170 **the decay factor and is compatible with distributional RL.**

2171 **a.** Model 1 with standard TD learning. Simulations were run using the receptors sensitivities
 2172 from the biophysical models and data-derived asymmetric scaling factors (see Methods 3.3).
 2173 The model's predictions capture the signatures in cue-evoked dopamine responses (left) and
 2174 value predictions (right) irrespective of the choice of the decay factor (β).

2175 **b.** Model 1 within the distributional RL framework (see Methods 3.3). The model's predictions
 2176 also capture the signatures in cue-evoked dopamine responses (left) and value predictions (right)
 2177 irrespective of the choice of the decay factor (β). Bottom row shows the distribution of value
 2178 predictors for each reward-predictive cue.

2179



2180

2181 **Extended Data Fig. 13 | Distributional reinforcement learning with D1 and D2 populations**
 2182 **(Model 1).**

2183 **a.** Schematic of the distributional RL model with D1 and D2 populations. The schematic
 2184 represents three different value predictors (pessimistic, neutral and optimistic from left to right)
 2185 with their respective P and N neurons. The level of optimism of each individual value predictor
 2186 is determined by the scaling factors of the individual dopamine RPE-evoked responses (α_i^+ , α_i^- ,
 2187 represented by the color in the colormap from purple to pink) and allows the model to encode
 2188 information about the distribution of rewards (bottom). The global level of ‘optimism’ or

2189 ‘pessimism’ of the agent is given by the re-scaling of the RPEs by the P and N receptors
2190 sensitivities in the model (ϕ^+ , ϕ^- , represented with the color saturation).

2191 **b.** Example of a Bernoulli distribution, equivalent to the reward distribution predicted by the
2192 50% cue.

2193 **c.** Distribution of expectiles learnt by the distributional RL model with D1 and D2 population for
2194 the reward distribution in b. The expectiles are sorted based on the asymmetric scaling factor of
2195 each individual dopamine neuron. Colormap represents the level of optimistic or pessimism of
2196 each agent.

2197 **d.** Samples from the decoded distributions for the set of expectiles in c. The probability density is
2198 bimodal, in accordance with the distribution in b. As the agents goes from pessimism to
2199 optimism, the probability density modes change in their relative magnitude.

2200