

1 An opponent striatal circuit for distributional reinforcement learning

2

3

4 Adam S. Lowet^{1,2,3}, Qiao Zheng^{1,4}, Melissa Meng^{1,2}, Sara Matias^{1,2}, Jan Drugowitsch^{1,4,*}, and
5 Naoshige Uchida^{1,2,*}

6 1. Center for Brain Science, Harvard University, Cambridge, MA, USA

7 2. Department of Molecular and Cellular Biology, Harvard University, Cambridge, MA, USA

8 3. Program in Neuroscience, Harvard University, Boston, MA, USA

9 4. Department of Neurobiology, Harvard Medical School, Boston, MA, USA

10 * Correspondence: uchida@mcb.harvard.edu (N.U.); jan_drugowitsch@hms.harvard.edu (J.D.)

11

12

13 **Abstract:** Machine learning research has achieved large performance gains on a wide range of
14 tasks by expanding the learning target from mean rewards to entire probability distributions of
15 rewards — an approach known as distributional reinforcement learning (RL)¹. The mesolimbic
16 dopamine system is thought to underlie RL in the mammalian brain by updating a representation
17 of mean value in the striatum^{2,3}, but little is known about whether, where, and how neurons in
18 this circuit encode information about higher-order moments of reward distributions⁴. To fill this
19 gap, we used high-density probes (Neuropixels) to acutely record striatal activity from well-
20 trained, water-restricted mice performing a classical conditioning task in which reward mean,
21 reward variance, and stimulus identity were independently manipulated. In contrast to traditional
22 RL accounts, we found robust evidence for abstract encoding of variance in the striatum.
23 Remarkably, chronic ablation of dopamine inputs disorganized these distributional
24 representations in the striatum without interfering with mean value coding. Two-photon calcium
25 imaging and optogenetics revealed that the two major classes of striatal medium spiny neurons
26 — D1 and D2 MSNs — contributed to this code by preferentially encoding the right and left tails
27 of the reward distribution, respectively. We synthesize these findings into a new model of the
28 striatum and mesolimbic dopamine that harnesses the opponency between D1 and D2 MSNs^{5–15}
29 to reap the computational benefits of distributional RL.

30 Main Text

31 Midbrain dopamine neurons and their primary target, the striatum, constitute an evolutionarily
32 ancient¹⁶ neural circuit that is critical for motivated behaviors^{17,18}. Computationally, dopamine
33 has long been thought to signal reward prediction error (RPE)^{2,19,20}, reminiscent of the teaching
34 signals used in many reinforcement learning (RL) algorithms²¹. Consistent with this idea,
35 dopamine is also known to modulate plasticity of certain corticostriatal synapses in roughly the
36 manner predicted by RL theory²², allowing neurons in the striatum to learn a representation of
37 average anticipated reward^{23–28}, often called “value”.

38 Despite the simplicity and popularity of this model, it leaves many aspects of the mesolimbic
39 circuit unexplained. For one, value representations reside not only in the striatum but throughout
40 the entire brain^{29–35}, and are enriched in neurons projecting to the striatum^{36,37}. Second, the
41 striatum is far from uniform, containing a variety of interneuron subtypes as well as D1 and D2
42 medium spiny neurons (MSNs), whose projection patterns differ³⁸ and whose plasticity is
43 modulated in opposite directions by dopamine^{22,39,40}. These differences at the receptor level
44 translate to opposite coding properties^{5,6} and effects on behavior^{7–15}, but interpreting their
45 distinct roles is complicated by the fact that these two populations often co-activate^{41–45}. Third,
46 dopamine activity is much more complex than a simple scalar RPE, varying both qualitatively
47 across dopamine projection systems^{46–49} and quantitatively within systems^{4,50,51}. Whether such
48 diversity is cause to revise RPE-based accounts of dopamine^{4,52,53} or discard them altogether^{54,55}
49 is currently the subject of intense debate.

50 In parallel to these questions about the neuronal representation of value, the striatum —
51 particularly the ventral striatum (VS, also referred to as the nucleus accumbens) — has long been
52 associated with decision-making under risk. VS lesions^{56–58} and dopaminergic drugs^{59,60} can both
53 impair risky decision-making, with some groups suggesting a particular role for VS D2
54 MSNs^{61,62}. Aberrant processing of risk, in turn, is thought to underlie many diseases associated
55 with these circuits, particularly addiction^{63–65}. Given this, it is perhaps surprising that, with a few
56 exceptions⁶⁶, conventional RL models of the basal ganglia ignore the role of risk, and most
57 theoretical investigations of uncertainty focus on sensory noise rather than intrinsic, irreducible
58 environmental stochasticity^{67–70}.

59 Borrowing from tremendous successes and popularity in machine learning^{71–73}, it was recently
60 proposed⁴ that the residual heterogeneity within RPE-coding dopamine neurons^{74,75} resembles
61 the predictions of a particular algorithm known as Expectile Distributional RL (EDRL)⁷⁶. This
62 algorithm not only unifies the learning of value and risk (and potentially higher-order moments
63 of reward distributions) within the same biologically-plausible framework but also provides
64 novel computational advantages — even with risk-neutral settings — related to representation
65 learning in deep neural networks^{77,78} and, potentially, directed exploration^{79–82}. However,

66 alternative accounts of the same dopamine data have since been put forward⁸³, including some
67 that question the very existence of a probabilistic value code^{84,85}.

68 Here, we provide the first direct evidence for distributional RL in the mammalian brain by
69 demonstrating that the striatum, and particularly VS, encodes not just mean value but also reward
70 variance. We combine our observations with well-established features of the anatomy and
71 physiology of the basal ganglia to construct a new computational model of reward distribution
72 learning in the striatum. The proposed model brings together diverse dopamine inputs⁴ and
73 asymmetric plasticity rules^{22,39,40} to enable a biological implementation of EDRL. Our model
74 makes several new experimental predictions about the representational geometry of the striatal
75 population code and its dependence on intact dopamine inputs, which we confirm using
76 Neuropixels recordings and dopamine lesions. Moreover, it suggests a way to unify the opponent
77 yet concurrent and non-redundant contributions of D1 and D2 MSNs to behavior via their coding
78 of the right and left tails of the reward distribution, respectively. We validate this view using cell
79 type-specific two-photon calcium imaging and optogenetic manipulations. Together, this study
80 improves our understanding of the computational principles underlying the brain's reward
81 circuitry and tightens the bonds between natural and artificial intelligence.

82 **A behavioral task to investigate distributional RL**

83 Representations of reward variance have been previously observed in a variety of cortical^{86–88}
84 and subcortical^{89–91} regions, but not in the striatum. To determine whether striatal neurons
85 encode reward variance while remaining agnostic to its representational format, we designed a
86 classical conditioning task in which water-restricted mice were trained to associate random odor
87 cues with probability distributions over stochastic reward magnitudes (Fig. 1a). Three different
88 probability distributions (Fig. 1b) were used: Nothing (100% chance of 0 μ L reward), Fixed
89 (100% chance of 4 μ L reward), and Variable (50/50% chance of 2/6 μ L reward). Fixed and
90 Variable distributions shared the same mean but had a different variance. Thus, distributional RL
91 predicts systematic differences in their underlying neural representations, whereas traditional RL
92 — assuming risk neutrality — does not. To ensure any such differences did not reflect
93 idiosyncratic odor preferences, two unique odors predicted each of the three distributions,
94 allowing us to compare odor representations both across- and within-distributions.

95 Crucially, while animals' anticipatory licking revealed a clear preference for Rewarded (Fixed
96 and Variable) over Unrewarded odors, it did not differ between the Fixed and Variable
97 distributions (Fig. 1c; here and elsewhere, we plot each mouse's mean across sessions as a
98 colored line for clarity, but the statistical tests disaggregate sessions using a Linear Mixed
99 Effects model with mouse-level random effects; see Methods). Additional behavioral data,
100 including face motion, whisking, pupil area, and running⁹², also did not support reliably
101 distinguishing Fixed from Variable trials (Fig. 1d and Extended Data Fig. 1a-b). The meager,
102 non-significant classification ability that may have existed was orthogonal to the regression

103 weight vector trained to predict value from all trial types (Extended Data Fig. 1d-e). This implies
104 that any ability to decode these trial types from neural data must be due to the associated
105 probability distributions and not to differential valuation or motor behavior.

106 **Striatum represents both mean and variance**

107 Next, we used high-density electrophysiological probes (Neuropixels) to acutely record activity
108 from across a broad swathe of the anterior striatum (Fig. 1e and Extended Data Fig. 2a; $N = 12$
109 mice, $n = 71$ sessions, 13,997 neurons). Consistent with prior work²³⁻²⁸, we found that both the
110 average firing rate of all neurons (Fig. 1f and Extended Data Fig. 2b) and the time course of trial
111 type-averaged activity projected onto the first principal component (PC; Fig. 1g) cleanly
112 separated Rewarded from Unrewarded odors. Furthermore, a substantial fraction of the activity
113 of individual neurons within our 1 s analysis window just before reward delivery correlated
114 significantly with expected reward, allowing us to reliably predict mean value from neural
115 (pseudo-) population activity across all striatal subregions (Extended Data Fig. 2c-e). Other
116 striatal neurons correlated significantly with reward prediction error during the reward period⁹³,
117 but these formed a smaller and mostly independent subset (Extended Data Fig. 2f-h).

118 However, not all neurons obeyed this simple pattern seen at the level of population averages.
119 Some single neurons consistently preferred Variable odors, while others — even when recorded
120 simultaneously — preferred Fixed (Fig. 2a). Such neurons fired similarly to *both instances* of the
121 Fixed and Variable odors, suggesting that they abstracted over odor-specific details to instead
122 encode information about variance — even as the population as a whole contained ample odor
123 information (Extended Data Fig. 3a-e).

124 Following observations of variance coding in other brain regions^{87,89}, we identified variance-
125 encoding neurons by linearly regressing single-neuron firing onto reward variance, after
126 regressing out the effect of mean reward. Unlike these prior studies, however, we surprisingly
127 found *fewer* striatal variance-encoding neurons than would be predicted from odor coding alone
128 (Extended Data Fig. 3f-h). Furthermore, in contrast to codes in which neural activity is construed
129 as representing samples from some probability distribution, across-trial Fano factors were the
130 same across trial types with different variances^{94,95} (Extended Data Fig. 4a-d). We therefore
131 adopted a different set of approaches to characterize distributional coding across the entire neural
132 population.

133 First, we projected each session's trial type-averaged firing rates in the 1 s window before reward
134 delivery ("Late Trace period") onto the first and second PCs (accounting for 72.9 ± 2.4 and 10.0
135 $\pm 1.0\%$ of the variance across trial types, respectively; mean \pm s.e.m. across mice; Fig. 2b). We
136 then measured the Euclidean distances in PC space along each dimension. As expected, trial
137 types with different mean rewards segregated out along PC 1 (Fig. 2c). More surprisingly,
138 though, Fixed and Variable odors separated out along PC 2, such that there was a greater
139 distance between across-distribution odor pairs than within-distribution odor pairs (Fig. 2d).

140 Second, to determine whether we could observe the same trends in native firing rate space, we
141 performed representational dissimilarity analysis (RDA) between the average population activity
142 vector for each of the rewarded trial types (Fig. 2e). Once again, the distance between across-
143 distribution pairs was greater, on average, than between within-distribution pairs (Fig. 2f). We
144 observed the same effects in the classification performance of single-trial linear classifiers
145 applied to pairs of rewarded trial types (Extended Data Fig. 5a-b) or applied to trial type groups
146 that either respected or violated their distribution identities (Extended Data Fig. 5c-d).
147 Distributional decoding was orthogonal to mean value coding (Extended Data Fig. 5e-h), stable
148 over time (Extended Data Fig. 5i-k), and strongest in the more ventral and lateral parts of the
149 striatum, particularly the lateral nucleus accumbens shell (lAcbSh; Extended Data Fig. 6). Lastly,
150 an artificial neural network-based decoder trained on single pseudo-trial population activity
151 successfully predicted complete reward distributions, even when its training and evaluation was
152 restricted to trials with the same mean, and generalized to unseen odors (Extended Data Fig. 7).

153 **Variance is encoded abstractly**

154 The preceding analyses show that the neural activities evoked by odors identifying the same
155 distribution are more similar to one another than to those evoked by odors identifying
156 distributions with the same mean but different variances. Let us now ask about the *relationship*
157 between Fixed and Variable odor representations. More specifically, is variance represented in
158 an “abstract format” — i.e., in a way that supports generalization to unseen situations⁹⁶? To find
159 out, we adapted two previously-defined metrics⁹⁶ to our task: parallelism score and cross-
160 condition generalization performance (CCGP). Both ask, in different ways, whether there is a
161 consistent direction in firing rate space that distinguishes low and high-variance cues (see
162 Methods).

163 The parallelism score is simply the average cosine similarity between the two difference vectors
164 pointing from Variable to Fixed population activity, one for each odor identifying the respective
165 distribution (Fig. 2g). Across sessions and mice, these difference vectors were significantly more
166 aligned than would be expected by chance (Fig. 2h). Similarly, a decoder trained on one Fixed
167 vs. Variable dichotomy and then tested on the held-out dichotomy achieved above-chance
168 CCGP, averaged across all four possible dichotomies (Fig. 2i-j). These analyses show that
169 variance is not just encoded arbitrarily, but in an abstract format.

170 **Using striatal opponency to implement distributional RL**

171 How might such an abstract representation be acquired? While there exist multiple theories for
172 how the brain might learn (factorized, that is, abstract) reward distributions^{71,72,83}, EDRL⁷⁶ is
173 especially promising because it requires only minimal modifications to existing, empirically
174 tested models of the basal ganglia⁴. EDRL proposes not just a single value predictor but an entire
175 family of predictors V_i , which learn at different rates, α_i^+ and α_i^- , for positive and negative RPEs,
176 respectively (Fig. 3a). “Optimistic” predictors have relatively high α_i^+ and will converge to

177 values above the distribution mean, while the opposite is true of “pessimistic” predictors. Each
178 predictor converges to a so-called “expectile” of the reward distribution, parameterized by $\tau_i =$
179 $\frac{\alpha_i^+}{\alpha_i^+ + \alpha_i^-}$ between 0 and 1. Expectiles generalize the mean ($\tau = 0.5$) just as quantiles generalize the
180 median, and collectively, they characterize the complete reward distribution⁹⁷ (Fig. 3b; see
181 Methods).

182 While EDRL has some appealing properties, it ignores the molecular and cellular diversity
183 within the striatum, most notably the presence of D1 and D2 MSNs³⁸. As an extension, we
184 propose reflected EDRL (REDRL) — so called because D2 MSN activity is simply the negative
185 of the corresponding value predictor, plus a constant offset to ensure non-negative activities (Fig.
186 3c). This simple modification does not merely lend EDRL additional biological plausibility;
187 rather, it demonstrates how the particular anatomy of the striatum can benefit distributional RL
188 computations while explaining a host of data regarding activity in the striatum and opponency
189 between D1 and D2 MSNs.

190 To implement REDRL in the striatum, we first require structured heterogeneity in dopamine
191 inputs, which can be modeled as piecewise linear response functions to reward size⁴ (Fig. 3d).
192 Since RPE is defined as actual minus predicted reward, the reward amount which elicits no
193 change in dopamine firing relative to baseline — the so-called “zero-crossing point”⁴ — is
194 equivalent to the learned value prediction for that neuron. Pessimistic dopamine neurons have
195 steeper slopes for rewards below their associated value prediction (α_i^-) and shallower slopes
196 above it (α_i^+), reflecting relatively low learning rates from positive RPEs. The converse is true of
197 optimistic dopamine neurons. Second, these diverse dopamine responses combine with opponent
198 plasticity rules in D1 and D2 MSNs, with D1 MSNs increasing synaptic weights more from
199 positive RPEs (β_m^+) and D2 MSNs increasing synaptic weights more from negative RPEs^{22,39,40}
200 (β_m^- ; Fig. 3e). Importantly, while asymmetric, these synaptic weight updates are not fully
201 dichotomous; D1 and D2 MSNs still learn slightly from dopamine changes in their non-preferred
202 directions⁶⁶, in line with the shallower but nonzero slope of D1 and D2 receptor occupancy
203 curves at baseline dopamine concentrations^{66,98,99}.

204 By composing these two functions, we get the complete REDRL model. The *opponency* of the
205 plasticity rule gives rise to opponent directions of value coding (Fig. 3f), with D1^{5,6,100,101} and D2
206 MSNs^{5,6,102} primarily correlating positively and negatively, respectively, with value and reward.
207 Meanwhile, its *asymmetric* nature has the effect of extremizing value predictors — D1 MSNs are
208 more optimistic, and D2 MSNs more pessimistic, than their individual dopamine inputs would
209 create on their own (Fig. 3g) — setting up a global bias for D1 and D2 MSNs to encode the right
210 and left tails of the value distribution, respectively, and shifting the zero-crossing points of the
211 coupled dopamine neurons up or down accordingly (Fig. 3d-g). Notably, this also predicts that
212 D1 MSNs will acquire positive associations faster than D2 MSNs, while D2 MSNs may be
213 preferentially involved in later discrimination or extinction^{39,40,103} (Fig. 3a).

214 Armed with such a model, we can ask whether the population activity predicted by REDRL
215 mirrors that observed in our striatal data. Strikingly, the top two PCs of the model predictors
216 closely resemble the projection of the data using principal component analysis (PCA; Fig. 3h-i).
217 Moreover, REDRL gives rise to a new prediction: Variable odors should be more distant from
218 Nothing odors along PC 1 than Fixed odors, a prediction that due to PCA's mean-centering is
219 independent of the D2 offset, and that we confirmed to hold true in our data (Fig. 3j-k).
220 Secondly, REDRL predicts the existence of substantial populations of neurons that correlate
221 either positively (D1) or negatively (D2) with expected reward across trials (Fig. 3l). We again
222 found this to be the case in our data (Fig. 3m), with a slight bias toward positive correlations,
223 perhaps reflecting the preponderance of D1 over D2 MSNs in the striatum¹⁰⁴⁻¹⁰⁶. Lastly, REDRL
224 predicts that the average firing rate should be slightly higher for Variable than for Fixed odors on
225 average, which we also observed (Fig. 3n-o). While other distributional RL formulations
226 predicted some of these effects, only REDRL and its close cousin, Reflected Quantile
227 Distributional RL (Extended Data Fig. 8a-m) predicted all of them. Thus, REDRL provides a
228 mechanistic account of distributional reinforcement learning which quantitatively matches the
229 structure of striatal representations.

230 **Dopamine is necessary for distributional RL**

231 If striatal representations are updated incrementally by dopamine RPEs as predicted by REDRL,
232 then eliminating dopamine prior to learning should disrupt these distributional representations
233 (Fig. 4a). To test this hypothesis, we injected the neurotoxin 6-hydroxydopamine (6-OHDA)
234 unilaterally into the lateral ventral striatum in naive mice, which resulted in local lesions of
235 dopamine neurons projecting to the injection site (Fig. 4b-c; Extended Data Fig. 9a). After
236 recovery, we trained the animals on the same task and then recorded neurons in both the control
237 and lesioned hemisphere ($N = 5$ mice, $n = 20$ sessions, 2,283 neurons from control; 19 sessions,
238 2,596 neurons from lesion). Unilateral lesions modestly impaired our ability to distinguish
239 Rewarded and Unrewarded odors based on behavioral predictors, but animals nonetheless
240 learned the task (Extended Data Fig. 9b-c).

241 Projecting striatal activity from each hemisphere independently into PC space suggested that
242 Fixed and Variable distributions were less well-separated in the lesioned hemisphere relative to
243 the control hemisphere (Fig. 4d). Indeed, when we quantified distances as before, we found
244 Nothing and Rewarded odors to be equally well separated along PC 1 for both hemispheres (Fig.
245 4e), but less-well separated along PC 2 in the lesioned hemisphere, with an associated smaller
246 difference in distances between across-distribution and within-distribution pairs (Fig. 4f).
247 Analogous effects were seen for parallelism score (Fig. 4g) and representational dissimilarity
248 (Fig. 4h), with stronger (and abstract) variance coding in the control relative to the lesioned
249 hemisphere. The persistence of mean value coding in the lesioned hemisphere may reflect the
250 inability of unilateral 6-OHDA to kill all dopamine neurons within the targeted hemisphere, the

251 interhemispheric broadcasting of mean value information once it reaches cortex^{31–37}, or, more
252 radically, the dispensability of dopamine for learning about mean value entirely.

253 In addition to supporting our mechanistic REDRL model, the selective disruption of variance
254 coding by 6-OHDA gives us an experimental tool with which to probe the function of
255 distributional RL in the brain. When paired with deep neural networks, distributional RL is
256 thought to boost system performance mainly by improving state representations^{1,4,78}. Due to
257 multiplexing of odor-specific representations alongside distribution information within the
258 striatum (Extended Data Fig. 3), it is possible to ask whether dopamine lesions also impair
259 striatal stimulus representations. We used multinomial logistic regression to decode odor identity
260 from neural activity during the 1 s window following odor onset. While we could decode odor
261 identity well above chance for both hemispheres, decoding performance was significantly higher
262 in the control than the lesioned hemisphere (Fig. 4i). The lesion caused a drop in decoding
263 performance across nearly all trial types, with the main driver being an increased confusion
264 between Fixed and Variable odors (Fig. 4j-k). These results are consistent with distributional RL
265 playing a similar role in shaping the representation of sensory inputs in artificial neural networks
266 and biological brains.

267 **Opponent contributions of D1 and D2 MSNs to REDRL**

268 To dissect the distinct contributions of D1 and D2 MSNs predicted by REDRL, we turned to
269 two-photon calcium imaging through implanted gradient refractive index (GRIN) lenses (Fig.
270 5a). We injected *AAV9-hSyn-FLEX-jGCaMP7s* virus¹⁰⁷ into the lAcbSh of *Drd1-Cre* ($N = 4$
271 mice, $n = 27$ sessions, 945 neurons) or *Adora2a-Cre* ($N = 4$ mice, $n = 38$ sessions, 1,106
272 neurons) transgenic mice¹⁰⁸, which drive expression in D1 and D2 MSNs, respectively (Fig. 5b).
273 Using this method, we were able to image up to ~50 neurons simultaneously per field of view
274 (31.6 ± 17.4 , mean \pm s.d. across sessions; Fig. 5c).

275 We observed different patterns of deconvolved Ca^{2+} activity across D1 and D2 populations^{100–102}
276 (Extended Data Fig. 10a-b). Many D1 MSNs were activated more to Rewarded than to
277 Unrewarded odors and outcomes, while the reverse was true, albeit less strongly, in D2 MSNs
278 (Fig. 5d-f). In addition to correlations across trials during the Late Trace period, we also
279 investigated differences between the Late Trace and Baseline periods, because REDRL predicts
280 increases in value predictions after Rewarded odor onset. Consistent with our model, significant
281 fractions of D1 and D2 MSNs increased and decreased their activities, respectively, more on
282 Rewarded than Unrewarded trials, although the pattern in D2 MSNs was again more
283 heterogeneous than in D1 MSNs (Extended Data Fig. 10c).

284 Intriguingly, we also found neurons which, like those we recorded using electrophysiology,
285 reliably distinguished between Fixed and Variable odors during the Late Trace period (Fig. 5g).
286 To test whether these trends were systematic, we performed the same analyses (RDA, CCGP and
287 PCA) separately on D1 and D2 MSNs, while pooling across all sessions and all mice to

288 compensate for the lower cell counts and higher variability of Ca^{2+} signals. Consistently across
289 disjoint subsets of trials in both D1 and D2 MSNs, across-distribution pairs were represented
290 more dissimilarly than within-distribution pairs (Fig. 5h), and variance was encoded in an
291 abstract format (Fig. 5i).

292 REDRL not only predicts that distributional coding should be present in both D1 and D2 MSNs
293 independently but also specifies the ways in which this coding should differ. For one, this
294 particular set of distributions should elicit higher variance across trial types for optimistic than
295 for pessimistic reward predictors on average — which is also true in the two-photon data for D1
296 and D2 MSNs, respectively (Fig. 5j). More impressively, when projecting optimistic ($\tau > 0.5$)
297 and pessimistic ($\tau < 0.5$) predictors into 2D PC space separately, we found that optimistic
298 predictors exhibited the same trend as the full complement of value predictors, with Variable and
299 Nothing odors further separated along PC 1 than Fixed and Nothing odors (Fig. 5k). However,
300 pessimistic predictors actually showed the opposite trend, with Variable and Nothing odors
301 closer together along PC 1 than Fixed and Nothing (Fig. 5l-m). Analogously, representational
302 dissimilarity was less for Variable and Nothing odors than Fixed and Nothing odors specifically
303 for pessimistic predictors; optimistic predictors did not differ (Fig. 5n). PCA projections (Fig.
304 5o-q) and RDA (Fig. 5r) for D1 and D2 MSNs mirrored these predictions precisely, revealing a
305 subtle distinction in distributional coding across MSN subtypes and confirming a novel
306 prediction of REDRL.

307 **Perturbing REDRL with optogenetics**

308 As a final test of REDRL, we sought to independently manipulate D1 and D2 MSNs while mice
309 performed a similar classical conditioning task. To do so, we expressed either the excitatory
310 opsin CoChR¹⁰⁹ ($N = 12$ mice, $n = 96$ sessions) or the inhibitory opsin GtACR1^{110,111} ($N = 13$
311 mice, $n = 92$ sessions) in D1 or D2 MSNs and implanted an optical fiber in lAcbSh¹¹² (Fig. 6a).
312 We then manipulated these neurons during the 2 s Trace Period after odor offset and quantified
313 licking during the last 1 s of this Trace Period, just prior to reward delivery (Fig. 6b).

314 To identify the REDRL model predictions for these manipulations, we clamped the simulated
315 values of inhibited and excited predictors respectively at 0 and 8 μL , the maximum reward size
316 we delivered in these experiments. We performed these simulated manipulations separately in
317 optimistic and in pessimistic neurons while letting the non-manipulated predictors retain their
318 original values. We then computed the animal's predicted value estimate as the mean across both
319 optimistic and pessimistic predictors (Fig. 6c-d). For comparison, we performed similar
320 manipulations on other distributional code types (Extended Data Fig. 11). We then took the
321 difference between the models' estimated values in Manipulation vs. No Manipulation trials for
322 each trial type (Fig. 6e; Extended Data Fig. 12) and compared it to the difference in anticipatory
323 licking. REDRL not only captured the main effects of “go” and “no-go” pathways¹¹³ but also
324 predicted precise patterns of licking across trial types, even for the same type of manipulation

325 (Fig. 6f). This could not be explained simply by ceiling effects, as the increase in licking was
326 sometimes greater for Rewarded than Unrewarded odors, as in the case of D2 inhibition.
327 Quantitative comparison between the data and various models confirmed that REDRL (and the
328 highly similar Reflected Quantile code) best fit the licking data (Fig. 6g).

329 Discussion

330 Here we have combined large-scale electrophysiology with cell-type specific recordings and
331 manipulations to develop the REDRL model of the basal ganglia. This model maintains the
332 normative algorithmic advantages of distributional RL¹ while lending itself to a biological
333 implementation that is consistent with the observed structure of dopamine population activity⁴
334 and dopamine-mediated plasticity rules^{22,39,40}. The most notable feature of REDRL is the distinct
335 role played by D1 and D2 MSNs, which specialize in the right and left tails of the reward
336 distribution, respectively. This bifurcated layout resembles other neural systems, such as
337 ON/OFF pathways in vision, and likely has similar benefits, such as efficient coding¹¹⁴, reduced
338 metabolic cost¹¹⁵, and flexibility¹¹⁶. For example, certain computations, such as expected value
339 estimation, would benefit from combining information from D1 and D2 MSNs, but others, such
340 as risk-sensitive behavior, might depend on one tail or the other, and thus primarily require
341 information from a single population. Furthermore, this architecture simplifies the problem of
342 connectivity: anatomically and/or genetically-defined subsets of dopamine neurons^{117,118} could
343 form independent closed loops with D2 (via ventral pallidum) and D1 MSNs, thereby helping to
344 keep separate pessimistic and optimistic RPE channels. These predictions should form the basis
345 of future anatomical investigations into the mesolimbic dopamine circuitry, as well as theories of
346 alternative architectures that might obviate this need¹¹⁹, which is shared by EDRL.

347 At the level of the striatum, REDRL helps unify previous approaches to understanding D1 and
348 D2 MSNs within a single, normative framework. While there have been previous hints that D1
349 and D2 MSNs are oppositely modulated by dopamine^{22,39,40} and oppositely correlated with
350 reward and expected value^{5,6,100–102}, this has generally been attributed to go/no-go or
351 approach/avoid pathways and modeled using a single value predictor^{3,66,113,120,121}. Here, we show
352 how, far from being a bug or redundancy in the RL architecture, such diversity could actually be
353 a feature, biasing convergence to optimistic or pessimistic value predictors. More speculatively,
354 it could explain why D1 and D2 MSNs often act in an opponent fashion without being inverses
355 of each other^{41–45}. The tendency for both pathways to activate prior to movement onset, for
356 example, would be predicted if such transition points coincide with increases in the predicted
357 variance of rewards (and thus the density on both the left and right tails).

358 REDRL also lends a new perspective to the coding of uncertainty in the brain. Typical treatments
359 of this topic focus on *perceptual* uncertainty, where the observer's role is to infer the distribution
360 of world states consistent with a pattern of neural activity⁷⁰. While the problem is generally
361 formulated as one of Bayesian inference⁶⁷, the associated uncertainty is frequently attributed to

362 noisy inputs rather than ones that are genuinely ambiguous (as in the case of the Necker cube¹²²).
363 In RL settings, in contrast, uncertainty generally arises from a combination of state ambiguity,
364 insufficient exploration, and intrinsic stochasticity¹²³, all of which complicate the problem of
365 learning from limited experience. Distributional RL excels in partitioning out this intrinsic
366 uncertainty from other sources, potentially allowing for improvements in state representation^{77,78},
367 exploration^{79–82}, value estimation¹²⁴, model-based learning¹²⁵, off-policy learning¹²⁶, and risk
368 sensitivity^{127–130}.

369 Many questions remain as to how the brain transforms high-dimensional reward distributions
370 into a single choice, but it is tempting to speculate that this process corresponds to the
371 dimensionality reduction that takes place throughout the various nuclei of the basal ganglia¹³¹,
372 ultimately collapsing onto a unitary value estimate in the mediodorsal thalamus that defines the
373 choice axis. Notably, such a “distributional critic” — centered here in the lAcbSh, a region
374 which receives RPE-like mesolimbic dopamine input^{46–50} — could integrate seamlessly into a
375 broader RL framework^{132–136}, with the dorsal striatum likely playing the role of the “actor” and
376 choosing actions in continuous, high-dimensional spaces¹³⁷. Modifications of the encoded reward
377 distribution, such as by dopaminergic drugs^{59,60}, or of the downstream basal ganglia circuit,
378 could then bias risky choice on rapid or developmental timescales^{61,138,139}. Various
379 psychopathologies — such as depression, in which patients learn more from losses than
380 gains^{140,141}, or addiction, in which patients systematically overweight the right tail of the reward
381 distribution¹⁴² — could similarly stem from the dysfunction of this core distributional RL
382 circuitry. Thus, REDRL can serve as a bridge between reinforcement learning, behavioral
383 economics, computational psychiatry, and systems neuroscience, demonstrating how the circuit
384 logic of the striatum can combine with vector-valued dopamine signals to realize the
385 computational benefits of distributional RL.

386 **Methods**

387 Experimental Procedures

388 *Mice*

389 A total of 46 adult C57BL/6J (Jackson Laboratory) male and female mice were used in these
390 experiments. Twelve wildtype animals (6 M, 6 F) were used for Neuropixels recordings, of
391 which five (2 M, 3 F) were also included in unilateral 6-OHDA experiments. For two-photon
392 imaging, four *Drd1-Cre* (B6.FVB(Cg)-Tg(Drd1-cre)EY262Gsat/Mmucd,
393 RRID:MMRRC_030989-UCD; 3 M, 1 F) and four *Adora2a-Cre* (B6.FVB(Cg)-Tg(Adora2a-
394 cre)KG139Gsat/Mmucd, RRID:MMRRC_036158-UCD; 1 M, 3 F) mice were used^{108,143,144}. For
395 optogenetic excitation, we used five *Drd1-Cre* (2 M, 3 F) and seven *Adora2a-Cre* (3 M, 4 F)
396 animals. For optogenetic inhibition, we crossed these lines with a Cre-dependent *GtACR1*
397 reporter mouse¹¹¹ (R26-CAG-LNL-GtACR1-ts-FRed-Kv2.1, RRID:IMSR_JAX:033089). Five
398 *Drd1-Cre;GtACR1* (2 M, 3 F) and eight *Adora2a-Cre;GtACR1* (4 M, 4 F) mice were used. All
399 transgenic mice used for experiments were backcrossed with C57BL/6J and heterozygous for the
400 relevant allele(s).

401 Animals were housed on a 12 hr dark/12 hr light cycle and performed the task at the same time
402 each day (\pm 1 hour), during the dark period. Ambient temperature was kept at $75 \pm 5^\circ\text{F}$, and
403 humidity was kept below 50%. Animals were group-housed (2–5 animals/cage) until surgery,
404 then individually housed throughout training and testing. All procedures were performed in
405 accordance with the National Institutes of Health Guide for the Care and Use of Laboratory
406 Animals and approved by the Harvard Institutional Animal Care and Use Committee (IACUC).

407 *Surgeries*

408 All surgeries were performed under aseptic conditions. Mice (> 8 weeks old) were anesthetized
409 with isoflurane (3.5% induction, followed by 1–2% maintenance at 1 L/min), and local
410 anesthetic (lidocaine, 2%) was administered subcutaneously at the incision site. Analgesia
411 (buprenorphine for pre-operative treatment, 0.1 mg/kg, intraperitoneal (i.p.); ketoprofen for post-
412 operative treatment, 5 mg/kg i.p.) was administered for two days after surgery. After leveling,
413 cleaning, and drying the skull, we affixed a custom-made titanium head plate to the skull with
414 adhesive cement²⁰ (C&B Metabond, Parkell).

415 For all injections, the solution (6-OHDA or virus) was backfilled into a pulled glass pipette
416 (Drummond, 5-000-1001-X), followed by mineral oil and a plunger. A small craniotomy (< 1
417 mm diameter) was made using a dental drill, and then the pipette assembly was mounted on the
418 stereotaxic holder, lowered to the desired coordinate, and injected slowly (~100 nL/min) to
419 minimize damage to the surrounding tissue (Narishige, MO-10). After each injection, we waited
420 at least 10 minutes to allow the solution to diffuse away from the pipette tip before slowly going

421 up to the next coordinate or retracting the pipette from the brain. Target coordinates (in mm) for
422 the lAcbSh were the same across experiments: AP 1.1 from bregma, ML 1.7, and DV 4.2 from
423 the pial surface.

424 *6-OHDA procedure*

425 To unilaterally ablate dopamine neurons projecting to lateral ventral striatum, we followed an
426 existing protocol^{52,145}. The following solution was injected (i.p.) into animals at 10 mg/kg prior
427 to surgery:

- 428 ● 14.25 mg desipramine (Sigma-Aldrich, D3900-1G)
- 429 ● 3.1 mg pargyline (Sigma-Aldrich, P8013-500MG)
- 430 ● 5 mL distilled water

431 Most animals (weighing ~25 g) received ~250 μ L of this solution, which was given to prevent
432 dopamine uptake in noradrenaline neurons and to increase the selectivity of uptake by dopamine
433 neurons. We additionally prepared a solution of 10 mg/mL 6-hydroxydopamine (6-OHDA;
434 Sigma-Aldrich, H116-5MG) and 0.2% ascorbic acid in saline (0.9% NaCl; Sigma-Aldrich,
435 PHR1008-2G). The ascorbic acid in this solution helps prevent 6-OHDA from breaking down.
436 The control hemisphere was either injected with vehicle ascorbic acid solution or uninjected; we
437 observed no differences between these groups and so combined them. To further prevent 6-
438 OHDA from breaking down, we kept the solution on ice, wrapped in aluminum foil, and used it
439 within three hours of preparation. If the solution turned brown during this time (indicating that 6-
440 OHDA had broken down), it was discarded and a fresh solution was made. 225 nL 6-OHDA (or
441 vehicle) was injected unilaterally into lAcbSh.

442 Surgeries occurred at least 1 week before the start of behavioral training. We lesioned nine
443 animals and included control hemisphere data for all of them in the main dataset. However, four
444 of these animals either died before we could record from the lesioned hemisphere or were not
445 correctly targeted for the lesion and/or recording, and so were excluded from the lesion dataset.

446 *Viruses*

447 To express constructs specifically in D1 or D2 MSNs, we injected viruses into *Drd1-Cre* and
448 *Adora2a-Cre* mice. For imaging experiments, we unilaterally injected 450 nL AAV9-hSyn-flex-
449 GCaMP7s ($\geq 1 \times 10^{13}$ vg/mL, Addgene)¹⁰⁷ into lAcbSh. For optogenetic activation experiments,
450 we bilaterally injected AAV9-hSyn-flex-CoChR-GFP (5.1×10^{12} vg/mL, UNC Vector Core,
451 NC)¹⁰⁹ at AP 1.1, ML ± 1.7 in 300 nL increments at four separate depths below the pial surface:
452 4.2, 3.4, 2.6, and 1.8.

453 *GRIN lens and fiber implantations*

454 Prior to GRIN lens surgery we injected animals i.p. with 50 μ L dexamethasone (2 mg/mL;

455 Vedco) to reduce inflammation. Before virus injection, a needle was mounted on the stereotaxic
456 holder, connected to light suction, and lowered to 3.4 mm below the pial surface to gently
457 aspirate away the overlying brain tissue. After virus injection, a singlet GRIN lens (0.5 NA, 0.6
458 mm diameter, 7.3 mm length, 0 – 200 μ m WD, 3/2 pitch, Inscopix, 1050-004597) was mounted
459 onto a stereotaxic cannula holder (Doric) and then slowly lowered over at least 30 minutes to its
460 target depth, 200 μ m above the injection site and 3.8 mm below the pial surface. Metabond was
461 used to secure the GRIN lens on all sides and allowed to dry completely before removing the
462 cannula holder and covering everything with another layer of Metabond mixed with charcoal
463 powder to block out light. Lastly, a plastic cap was attached with Kwik-Cast (World Precision
464 Instruments) to protect the lens from damage.

465 For optogenetic manipulation, we bilaterally implanted tapered fibers (0.66 NA, 200 μ m
466 diameter, 3 mm emitting length, 5 mm implant length; Optogenix) in the lAcSh after virus
467 injection, at a depth of 4 mm. Each fiber was secured using Metabond and then protected with a
468 fitted cap.

469 *Behavior setup and tasks*

470 Behavioral events were controlled (and licking was monitored) using custom-written software in
471 MATLAB (Mathworks, Natick, MA) and the Bpod library (Sanworks, Rochester, NY)
472 interfacing with the Bpod state machine (Sanworks, 1024 and 1027), valve module (Sanworks,
473 1015), and port interface board (Sanworks, 1020)/water valve (Lee Company, LHDA1233115H)
474 assembly. Odors were delivered using a custom olfactometer¹⁴⁶, which directed air through one
475 of eight solenoid valves (Lee Company, LHDA1221111H) mounted on a manifold (Lee
476 Company, LFMX0510528B). Each odor was dissolved in mineral oil at 10% dilution, and 30 μ L
477 of diluted odor solution was applied to a syringe filter (2.7 μ m pore, 13 mm diameter; Whatman,
478 6823-1327). Wall air was passed through a hydrocarbon filter (Agilent Technologies, HT200-4)
479 and split into a 100 mL/min odor stream and 900 mL/min carrier stream using analog flowmeters
480 (Cole-Parmer, MFLX32460-40 and MFLX32460-42), which were recombined at the odor
481 manifold before being delivered to the animal's nose. Licking was monitored using an infrared
482 emitter-photodiode pair positioned just in front of the plastic lick spout, positioned at the
483 animal's mouth.

484 Animals used for Neuropixels recording and 2-photon imaging were conditioned with six
485 different neutral odors, chosen at random from these seven: isoamyl acetate, *p*-cymene, ethyl
486 butyrate, (*S*)-(+)-carvone, (\pm)-citronellal, α -ionone, and L-fenchone. Optogenetic manipulation
487 animals used only the first three. In all experiments, the mapping between physical odor and
488 conceptual trial type was randomized across mice. Each trial began with a 1 s odor presentation,
489 followed by 2 s trace period and then reward delivery. There was a minimum of 4.6 s before the
490 next trial (4.1 s for optogenetic manipulation animals), plus a variable ITI drawn from a
491 truncated exponential distribution with a mean of 2 s, minimum of 0.1 s, and maximum of 10 s.

492 For 2-photon imaging experiments, this was extended to a mean of 10.5 s, minimum of 6.5 s, and
493 maximum of 18.5 s to account for the slower kinetics of the calcium indicator relative to
494 electrophysiology.

495 The recording task consisted of three different reward distributions, Nothing, Fixed, and Variable
496 (Fig. 1b). Each distribution was then paired with two unique odors, for a total of six odors. The
497 distributions were as follows:

- 498 ● Nothing: 100% chance of 0 μL water
- 499 ● Fixed: 100% chance of 4 μL water
- 500 ● Variable: 50% chance of 2 μL water; 50% chance of 6 μL water

501 The task used for optogenetic manipulation was simplified in two ways. First, we used only one
502 odor per distribution, for a total of three odors. Second, we modified the Variable distribution to
503 be 50/50% between 0 and 8 μL , because our model predicted that increasing the variance would
504 lead to a greater behavioral difference between Fixed and Variable odors.

505 *Behavior training*

506 Water restriction began no earlier than 5 days after recovery from surgery. Animals' condition
507 was monitored daily to ensure that mice did not dip below 85% of their free-drinking body
508 weight, including supplementing with additional water after the task to bring their total daily
509 intake to ~ 1.2 mL. Over the course of three successive habituation days, mice were (1) handled
510 gently for several minutes in their home cage, (2) permitted to freely roam around the platform in
511 the behavior rig to collect water and then (3) head-fixed while receiving frequent (inter-reward
512 interval 4-5 s) 6 μL water rewards.

513 The optogenetic manipulation task proceeded in only one phase, with up to 110 Nothing, 110
514 Fixed, and 114 Variable trials, randomly interleaved. By contrast, training for the recording task
515 took place in three phases, each with a maximum of 300 trials.

- 516 ● In Phase 1, mice experienced both Nothing odors and both Fixed odors with equal
517 probabilities
- 518 ● In Phase 2, mice experienced all six odors, but with the Variable odors 5.5x more
519 frequent than the others
- 520 ● In Phase 3, mice experienced all six odors at the final ratio of 4:4:7
521 (Nothing:Fixed:Variable), to increase the statistical power for analyzing responses to
522 different reward sizes

523 On recording days, animals experienced a maximum of 20 additional Unexpected reward trials,
524 in which 4 μL of water was delivered without being preceded by an odor cue. All trials were
525 randomly interleaved in all phases.

526 For both tasks, animals completed at least 150 trials per day, and almost always more than 250.
527 The experiment might be terminated early by the experimenter if the animals stopped licking in
528 anticipation (or consumption) of the rewards due to satiety. A behavior session was considered
529 “significant” if the lick rate during the last half second prior to reward delivery was significantly
530 different between Rewarded (Fixed and Variable) and Unrewarded (Nothing) odors (Mann-
531 Whitney U test, $\alpha = 0.05$) and the effect size was at least 0.75 licks/s. Animals were advanced to
532 the next phase, or to habituation for recording/manipulation, after at least two consecutive days
533 with significant behavior. On recording/manipulation days, only significant behavior sessions
534 were included for neural or behavioral analysis.

535 *Neuropixels recordings*

536 The day before recording, animals were habituated to the recording setup by covering their heads
537 with a plastic sheet to block their view of the probe and manipulator. We then turned on the
538 lamp, ran the brushed motor controller (Thorlabs, KDC101 and Z825B) up and down several
539 times, tapped on the skull several times with fine forceps, and left the animal head-fixed for at
540 least 30 mins before beginning the behavioral protocol. If necessary, we repeated this habituation
541 protocol every day until the animal’s behavior was significant (see “Behavioral training” above).
542 After this, we anesthetized the animal to make a small craniotomy, which was then covered with
543 Kwik-Cast. The craniotomy was guided by fiducial marks made at the target sites for probe
544 insertion during headplate implantation using a fine-tipped pen. Target coordinates included: AP
545 0.9, ML 1.7 (lAcbSh); AP 1.1 ML 1.4 (nucleus accumbens core); and AP 1.4, ML 0.6 (medial
546 accumbens shell, mAcbSh). For the first craniotomy, a ground pin was inserted into the posterior
547 cortex and a custom-made plastic recording chamber was fixed to the top of the headplate, both
548 using five-minute epoxy (Devcon).

549 The next day, we head-fixed the mouse, covered its head as before, removed the Kwik-Cast, and
550 flushed the craniotomy with saline. For the first recording in each craniotomy, we coated the
551 probe in lipophilic dye at 10 mg/mL. DiI (1,1'-dioctadecyl-3,3,3',3'-tetramethylindocarbocyanine
552 perchlorate, Sigma-Aldrich, 42364-100MG) and DiD (1,1'-dioctadecyl-3,3,3',3'-
553 tetramethylindocarbocyanine, 4-chlorobenzenesulfonate, Biotium, 60014-10mg) were
554 dissolved in 100% ethanol (Koptec, V1001), and DiO (3,3'-dioctadecyloxycarbocyanine
555 perchlorate, ThermoFisher, D275) was dissolved in 100% *N,N*-dimethylformamide (Sigma-
556 Aldrich, D4254). The coated Neuropixels 1.0¹⁴⁷ or four-shank Neuropixels 2.0¹⁴⁸ probe was then
557 mounted on the manipulator, and connected to the ground pin via a wire soldered onto the
558 reference pad and shorted to ground. In the event the external reference was unstable, we used tip
559 referencing instead. All recordings were performed in SpikeGLX software
560 (<https://github.com/billkarsh/SpikeGLX>) with sampling rate = 30 kHz, LFP gain = 250, and AP
561 gain = 500, and we analyzed only the AP channel (which was high-pass filtered in hardware with
562 a cutoff frequency of 300 Hz).

563 We inserted the probe into the brain at 9 $\mu\text{m/s}$ before slowing to 2 $\mu\text{m/s}$ when we were 500 μm
564 above the target depth. We stopped insertion when we saw ventral pallidal activity, characterized
565 by large-amplitude, high-frequency spikes, on the first 40 channels or so (or 5 channels for
566 Neuropixels 2.0). This point was usually reached around 5.2 mm below the visually-identified
567 pial surface. After reaching the target depth, the probe was allowed to settle for 30 minutes prior
568 to starting the experiment and Neuropixels recording. Behavioral and neural recordings were
569 synchronized using a TTL pulse sent from the Bpod to the PXIe acquisition module SMA input
570 at the start of every trial. After the experiment, the probe was retracted at 9 $\mu\text{m/s}$ and the
571 craniotomy was re-sealed with Kwik-Cast. Neuropixels data were spike sorted offline with
572 Kilosort 3¹⁴⁹ with default parameters, followed by manual curation in Phy
573 (<https://github.com/cortex-lab/phy>).

574 *Two-photon imaging*

575 Imaging data were acquired using a custom-built two-photon microscope. A resonant scanning
576 mirror and galvanometric mirror (Cambridge Technology, CRS 8 KHz and 6210H) separated by
577 a scan lens-based relay on the scan head (Thorlabs, MM201) allowed fast scanning through a
578 dichroic beamsplitter (757 nm long-pass, Semrock) and 20x/0.5 NA air immersion objective lens
579 (Nikon, Plan Fluor). Green and red emission light were separated by a dichroic beamsplitter (568
580 nm long-pass, Semrock) and bandpass filters (525/50 and 641/75 nm, Semrock) and collected by
581 GaAsP photomultiplier tubes (Hamamatsu, H7422PA-40) coupled to transimpedance amplifiers
582 (Thorlabs, TIA60). A diode-pumped, mode-locked Ti:sapphire laser (Spectra-Physics) delivered
583 excitation light at 920 nm with an average power of ~ 60 mW at the top face of the GRIN lens¹⁵⁰,
584 modulated by a Pockels cell (Conoptics, 350-80). The microscope was controlled by ScanImage
585 (Version 4; Vidrio Technologies). The behavior platform was mounted on an XYZ translation
586 stage (Thorlabs, LTS150 and MLJ050) to position the mouse under the objective, and the top
587 face of the GRIN lens was first located using a 470 nm LED (Thorlabs, M470L2).

588 Due to the limited axial resolution of the implanted GRIN lens, we acquired only a single
589 imaging plane at 15.2 Hz unidirectionally with 1.4x digital zoom and a resolution of 512 x 512
590 pixels (~ 1 $\mu\text{m/pixel}$ isotropic). Imaging was either continuous or triggered 2.6 s before
591 odor/unexpected reward onset, depending on the session. Bleaching of GCaMP7s was negligible
592 over this time. TTL pulses were sent from the microscope to Bpod to synchronize imaging and
593 behavioral data. Imaging typically began ~ 4 weeks after GRIN implantation, to allow sufficient
594 time for the virus to express and for inflammation to clear.

595 *Two-photon pre-processing*

596 We used the Suite2p toolbox¹⁵¹ (version 0.10.3) to register frames, detect cells, extract Ca^{2+}
597 signals, and deconvolve these traces. We used parameter values of $\text{tau}=2.0$ (to approximately
598 match the decay constant of GCaMP7f¹⁰⁷), $\text{sparse_mode}=\text{False}$, $\text{diameter}=20$, $\text{high_pass}=75$,
599 $\text{neucoeff}=0.58$; fs was set to the measured frame rate for that session (~ 15.2 Hz), and all other

600 parameters were set to their defaults. Briefly, non-rigid motion correction was used in blocks of
601 128 x 128 pixels to register all frames to a common reference image using phase correlation. Cell
602 detection consisted of finding and smoothing spatial PCs and then extending ROIs spatially
603 around the peaks in these PCs. Next, Ca²⁺ traces were extracted from each ROI after discarding
604 any pixels belonging to multiple ROIs. Finally, neuropil contamination and deconvolved spikes
605 were estimated in a single step from Ca²⁺ fluorescence in each ROI using the OASIS
606 algorithm¹⁵² with a non-negativity constraint. This deconvolved activity was used for all
607 subsequent analysis. ROIs were manually curated on the basis of anatomical and functional
608 criteria using the Suite2p GUI to exclude neuropil and ROIs with few or ill-formed transients.

609 *Face and body imaging*

610 In addition to the lick port, we monitored behavior using two cameras at 30 Hz, one pointed at
611 the face (PointGrey, FL3-U3-13Y3M) and one pointed at the body (PointGrey, CM3-U3-13S2C)
612 under both visible and infrared LED illumination. Cameras were synchronized from Bpod once
613 per trial using GPIO inputs, and data were written to disk via Bonsai¹⁵³. Behavioral features were
614 extracted using custom code alongside Facemap⁹² (version 0.2.0). Face motion energy was
615 computed as the absolute value of the difference between consecutive frames and summed across
616 all pixels to yield the “whisking” signal. In addition, we performed singular value decomposition
617 (SVD) on the motion energy video (in chunks, following ref.⁹²) and projected the movie onto the
618 top 50 components to obtain their activity patterns over time. Pupil area was estimated simply as
619 the mean (inverse) pixel value within a mask, after interpolating over blink events. Running was
620 computed using the phase correlation of the cropped body video, to take into account limb and
621 tail movements.

622 *Optogenetic manipulation*

623 473 nm laser light (Laserglow Technologies, LRS-0473-GFM-00100-03) was delivered to the
624 implanted tapered fibers using a custom-built rig (modeled after refs.^{154,155}) coupled to a high-
625 performance patch cord (0.66 NA, Plexon, OPT/PC-FC-LCF-200/230-HP-2.2L KIT). Briefly,
626 light was split into two identical paths using a 50/50 beamsplitter cube (Thorlabs, CCM1-
627 BS013). Each path was then focused onto a galvanometric mirror (Novanta 6210K) and re-
628 collimated using an achromatic doublet (Thorlabs, AC508-100-A-ML), before being focused
629 onto the back of the patch cord using an aspheric condenser lens (Thorlabs, ACL50832U). This
630 setup allowed us to modulate the angle at which light entered the patch cord, and thus the
631 distance at which it exited the tapered fiber. We delivered light at two different angles (three in
632 some experiments), but here we analyze only ventral manipulation trials, in which the incident
633 angle of light was ~0°, light exited near the tip of the fiber, and coupling between the patch cord
634 and fiber was approximately 50%¹⁵⁴.

635 The laser output (and the angle of the galvanometric mirrors) was controlled by Bpod via
636 PulsePal¹⁵⁶ (Version 2; Sanworks, 1102). Stimulation was delivered bilaterally during the two

637 second-long trace period, immediately prior to reward. For CoChR excitation experiments, we
638 used 10 ms pulses at 20 Hz with an output power at the tapered fiber of 100 μ W. For GtACR1
639 inhibition, we used a constant, 1 mW pulse for the full 2 seconds. In both cases, stimulation was
640 delivered on 45.5% of trials, uniformly at random across manipulation locations and trial types.

641 *Histology and immunohistochemistry*

642 Mice were deeply anesthetized with ketamine/dexmedetomidine (80/1.1 mg/kg) and then
643 transcardially perfused using 4% paraformaldehyde. The brains were sliced at 100 μ m into
644 coronal sections using a vibratome (Leica) and stored in PBS. If performing immunostaining,
645 slice thickness was 75 μ m. These slices were then permeabilized with 0.5% triton X-100,
646 blocked with 10% FBS, and stained with rabbit anti-tyrosine hydroxylase antibody (TH; AB152,
647 EMD Millipore, RRID: AB_390204) at 1:750 dilution at 4°C for 24 hours to reveal dopamine
648 axons in the striatum. Next, slices were stained with fluorescent secondary antibodies (Alexa
649 Fluor 488 goat anti-rabbit secondary antibody, A-11008, Invitrogen, RRID: AB_143165) and
650 DAPI at 1:500 dilution at 4°C for 24 hours. Slices were then mounted on glass slides
651 (VECTASHIELD antifade mounting medium, H-1000, or with DAPI for non-stained slices, H-
652 1200, Vector Laboratories) and imaged using Zeiss Axio Scan Z1 slide scanner fluorescence
653 microscope. We visually verified the placement of all GRIN lenses and fibers to be within the
654 lAcbSh.

655 Data Analysis

656 *Atlas registration*

657 For electrophysiology experiments, we registered slices to the Allen Mouse Brain Atlas with
658 SHARP-Track¹⁵⁷ and used it to trace dyed probe trajectories in the AP and ML directions as well
659 as visualize the registered trajectories as a coronal stack. We also used this registration to define
660 the unique DV extent of each mouse's lateral ventral striatal 6-OHDA lesion, and we considered
661 only neurons that fell within this range to have been lesioned. To more accurately ascertain the
662 depth of recordings, we used the International Brain Lab's Ephys Atlas GUI
663 (<https://github.com/int-brain-lab/iblapps/tree/master/atlaselectrophysiology>), focusing on the
664 boundary between the ventral pallidum and nucleus accumbens due to the abrupt change in
665 electrophysiological characteristics at this interface. When necessary, we also adopted their
666 convention that in Allen Common Coordinate Framework¹⁵⁸ (CCF) coordinates, bregma = 5400
667 AP, 332 DV, and 5739 ML. For plotting probe trajectories in 3D, we used the Brainrender
668 library¹⁵⁹.

669 For more fine-grained analysis of subregions, we used the Kim Lab atlas¹⁶⁰ accessed through the
670 BrainGlobe Atlas API¹⁶¹. This atlas applies the Franklin and Paxinos¹⁶² labels to the Allen
671 CCF¹⁵⁸, with additional striatal subregions defined by Hintiryan et al.¹⁶³. For some subregions,
672 the parcellation was finer than we needed, so we pooled subregions as follows:

- 673 ● Olfactory tubercle (OT): Tu1; Tu2; Tu3
- 674 ● Ventral pallidum (VP): VP
- 675 ● Medial nucleus accumbens shell (mAcbSh): AcbSh
- 676 ● Lateral nucleus accumbens shell (lAcbSh): lAcbSh; CB; IPACL
- 677 ● Nucleus accumbens core (core): AcbC
- 678 ● Ventromedial striatum (VMS): CPr, imv; CPi, vm, vm; CPi, vm, v; CPi, vm, cvm
- 679 ● Ventrolateral striatum (VLS): CPr, l, vm; CPi, vl, imv; CPi, vl, v; CPi, vl, vt; CPi, vl, cvl
- 680 ● Dorsomedial striatum (DMS): CPr, m; CPr, imd; CPi, dm, dl; CPi, dm, im; CPi, dm, cd;
- 681 CPi, dm, dt
- 682 ● Dorsolateral striatum (DLS): CPr, l, ls; CPi, dl, d; CPi, dl, imd

683 *Unit inclusion criteria*

684 To be included for analysis, units from Neuropixels recordings had to have a minimum firing
685 rate of 0.1 Hz and to have been stable, defined as a coefficient of variation of firing rate
686 (computed in 10 equally-sized, contiguous, disjoint blocks during the session) less than 1. 13,997
687 single units survived these inclusion criteria in the main dataset. In the lesion dataset, we
688 additionally filtered neurons by their DV position: only those that fell within the DV range of the
689 lesion were included in the matched control dataset for that mouse. Of the 9,081 neurons that
690 survived the electrophysiological criteria, 4,879 were in the correct anatomical location, of which
691 2,283 came from the control and 2,596 came from the lesioned hemisphere.

692 *Putative cell type identification*

693 We assigned units to putative cell types using previously-established criteria¹⁶⁴. Briefly, to be
694 considered MSNs, units were required to have broad waveforms (Kilosort template trough-to-
695 peak waveform duration > 400 μ s) and post-spike suppression \leq 40 ms. For the latter, we used
696 the autocorrelation function with a bin width of 1 ms. Post-spike suppression was quantified as
697 the duration for which the autocorrelation function was less than its average during lags between
698 600–900 ms.

699 *Statistical software*

700 All statistical analysis, except where explicitly stated, was performed in Python using the NumPy
701 (v. 1.22.3), SciPy (v. 1.7.3), pandas (v. 1.1.4), scikit-learn (v. 1.0.2), statsmodels (v. 0.14.0),
702 Matplotlib (v. 3.5.1), and seaborn (v. 0.12.2) packages^{165–171}. If not otherwise specified,
703 statistical tests used Linear Mixed Effects models (LMEs) with a random intercept for each
704 mouse, and, if applicable, a random slope for each mouse as a function of grouping (e.g. Across-
705 vs. Within-distribution), implemented in statsmodels. All reported *p*-values are two-tailed.

706 *Units of analysis*

707 For the behavior, control and manipulation datasets (Figs. 1, 2, 3, and 6), each observation was
708 an individual session — that is, we used simultaneously-recorded neurons and behavior and
709 computed effects (PCA, RDA, parallelism score, classification) on a session-by-session basis.
710 However, given the limited spatial extent of our lesion and our lower number of simultaneously-
711 recorded neurons, for the lesion dataset (Fig. 4) we used pseudo-populations. More specifically,
712 we created pseudo-populations by splitting the dataset into disjoint sets of *trials*¹⁷², which were
713 stitched across sessions, but not across animals. Within each session, we used simultaneously-
714 recorded trials across neurons to preserve noise correlations where possible. For these LMEs
715 then, pseudo-populations provided the observations and mouse was again the grouping variable.
716 The same procedure was used for all subregion-specific analyses (Extended Data Figs. 2d, 3e,
717 6a-d) and ANN-based decoding (Extended Data Fig. 7a-d) due to the lower number of
718 simultaneously-recorded neurons available for these analyses.

719 For the imaging dataset (Fig. 5) and ANN-based transfer (Extended Data Fig. 7e-f), we did not
720 have enough neurons in all animals to assess distributional coding. We therefore pooled neurons
721 not only across sessions but also across animals within genotype. Pseudo-populations were
722 otherwise constructed exactly as in the lesion case. To be consistent with the parametric nature of
723 LMEs while recognizing that observations were no longer specific to individual mice, we used
724 one sample *t*-tests to assess statistical significance relative to chance levels and LMEs (with just
725 one observation per group) to assess differences between groupings.

726 The only exception to these choices was when computing the fraction of cells significantly
727 encoding each variable of interest (mean, reward, RPE, etc.), or their conjunction. In this case,
728 we always pooled across-sessions within-mouse, since we were computing a single fraction, and
729 used paired samples *t*-tests between data and shuffled fractions (or actual combined cells versus a
730 prediction assuming independence).

731 *Time periods for analysis*

732 In general, we analyzed behavioral and neural data during the Late Trace period, 1–0 s before
733 reward delivery. However, for odor decoding, we used the Odor period (0–1 s after odor onset),
734 and reward or RPE we used the Outcome period (0–1 s after reward delivery). Neural and
735 behavioral data were averaged within these 1 s periods before analysis, with the exception of
736 plots of classification or regression time courses, in which averages within non-overlapping 250
737 ms bins were used.

738 *Visualization of neural time courses*

739 For smoothed plots of neural time courses (Figs. 1f, g; 2a; 5d, g; Extended Data Fig. 2b, 10a-b),
740 we smoothed neural activity (spike trains or deconvolved activity traces) with a Gaussian kernel
741 (s.d. 100 ms) before plotting or reducing dimensionality. Z-scored firing rates were computed
742 using the mean and standard deviation of this smoothed trace. PCA time courses (Fig. 1g) were

743 extracted by computing the average normalized, smoothed firing rate for each trial type and
744 concatenating these into a 2D matrix of shape $N \times (T \times 6)$, where N is the number of neurons, T is
745 the number of time points per trial, and 6 corresponds to the six possible odors. PCA was then
746 performed and the time courses were reconstructed separately for each of the six odors. All other
747 analyses used unsmoothed data so as to not be contaminated by later time points.

748 *Principal component analysis and representational dissimilarity analysis*

749 For two-dimensional PC plots, normalized activity during the Late Trace period was averaged
750 across trials within a given type to produce a matrix of shape $N \times 6$. We then applied PCA to
751 reduce this matrix to shape 2×6 , having retained only the top 2 PCs. Results were qualitatively
752 identical when using all neurons or only putative MSNs for the main dataset (Fig. 2). We report
753 Euclidean distances between projected trial types, measured separately along each PC. RDA was
754 similar, except that we computed cosine distances in the native (pseudo-)population normalized
755 firing rate space, rather than a lower-dimensional projection.

756 *Parallelism score*

757 Following ref.⁹⁶, we computed the normalized mean firing rate in response to each of the Fixed
758 and Variable odors. There are two possible ways to pair up these four odors: (1) Fixed 1 vs.
759 Variable 1 and Fixed 2 vs. Variable 2, or (2) Fixed 1 vs. Variable 2 and Fixed 2 vs. Variable 1. In
760 both cases, we can compute difference vectors pointing from Variable to Fixed (Fig. 2g) and
761 then take the cosine similarity between them. The parallelism score we report is simply this
762 cosine similarity, averaged over the two possible divisions. Note that in the case of isotropic
763 noise, the vectors that we define are equivalent to those defined by a maximum-margin linear
764 classifier between the two conditions. However, high parallelism score does not necessarily
765 imply high cross-condition generalization performance (CCGP) — for example, if the test
766 conditions are much closer together than the training conditions, the noise is high and/or
767 anisotropic, or the coding directions for different variables are not orthogonal (e.g. arranged as a
768 parallelogram rather than a rectangle).

769 *Classification*

770 For both behavioral and neural binary classification, we used a support vector classifier (SVC)
771 with a linear kernel, hinge loss function, L2 penalty, balanced accuracy scoring across classes,
772 and regularization parameter 5×10^{-3} , implemented in scikit-learn. The linear kernel allows for
773 easy interpretation of the learned weights. Input data (unnormalized spike counts, lick counts, or
774 mean Facemap predictors) were transformed using StandardScaler (computed on training data)
775 before being fed to the classifier.

776 We ran five different classification analyses: CCGP⁹⁶, pairwise decoding, congruency, mean, and
777 odor, as described in the Main Text and figure legends. Across-distribution and within-

778 distribution results were just the average over the relevant dichotomies (e.g. the four possible
779 ways to set up CCGP). For all simultaneous decoding analyses except for CCGP, five cross-
780 validation folds were used, and reported classification accuracy was the average over these five
781 folds. For CCGP, cross-validation was unnecessary because training and test sets were fully
782 disjoint already. Similarly, for pseudo-population based decoding (Figs. 4–5), 5 training sets and
783 1 disjoint test set were used in all cases. For six-way odor classification, we used multinomial
784 logistic regression rather than SVC, again with a regularization parameter of 5×10^{-3} and
785 balanced accuracy scoring across classes.

786 Cross-temporal decoding (Extended Data Fig. 3d, 5h-j) settings were identical to the above. For
787 the odor, pairwise, and congruency analyses, we ensured that the same trial never appeared in
788 both the training and testing sets, despite the different time windows used, to avoid leakage due
789 to temporal autocorrelation. For CCGP, train and test trials were always different, so this was not
790 a concern.

791 *Cosine similarity to classification boundary*

792 Both linear classification and regression find a high-dimensional weight vector in neural state
793 space; computing the cosine similarity between these vectors can identify whether two analyses
794 are honing in on the same or different features. For each session, in addition to performing
795 classification as described above, we regressed input data (unnormalized spike counts, lick
796 counts, or mean Facemap predictors) during the same time period against per-trial mean or
797 variance (using StandardScaler followed by RidgeCV with default scikit-learn parameters). Note
798 that the regression uses all six trial types, while the classification is limited to looking at only two
799 (pairwise or CCGP) or four (congruency or mean) odors at a time. We then took the weights
800 learned by each regression and computed the cosine similarity with the classification weights
801 (separately for each of the five classification cross-validation folds for non-CCGP decoders; each
802 session was summarized as the average of these five measurements). We report the results of an
803 LME testing either the difference from a chance value of 0, indicating orthogonality (CCGP), or
804 the difference between the absolute cosine similarities for across- and within-distribution
805 decoders (pairwise and congruency; Extended Data Fig. 5f-g).

806 *Distribution-coding subpopulation*

807 To identify neurons that contributed significantly to distribution decoding, we extracted the
808 coefficients from each session's CCGP, pairwise, and congruency decoders and averaged them
809 across dichotomies (and across cross-validation folds if necessary). For the pairwise and
810 congruency analyses, we additionally took the difference between Across- and Within-
811 distribution coefficients. For each quantile level (computed on each set of coefficients
812 individually for each mouse and each decoder), we then calculated the fraction of neurons above
813 this quantile level for all three decoders compared to null decoders in which trial types had been
814 shuffled before being run through the decoder. We chose a cutoff such that only 2.5% of these

815 cells from the null decoders survived; for the actual data, this corresponded to 1,600 significant
816 distribution-coding neurons, or 11.43% of the total. We refer to these neurons as the
817 “distribution-coding subpopulation” (Extended Data Fig. 6e-f, 7).

818 *Percentage of significant cells*

819 To compute correlations with different variables of interest, we calculated the trial-wise Pearson
820 correlation between unsmoothed activity in a given bin and the value of the variable of interest
821 on that trial. We then did the same thing, except that for each neuron independently we shuffled
822 the mappings between odor and distribution. For example, when considering correlations with
823 mean value, a Fixed 1 trial would correspond to a mean of 4 (μL). If upon shuffling, Fixed 1
824 odors were mapped to Nothing 2, then the corresponding mean in the shuffled dataset would be
825 0. Percentages of cells significantly correlating with variables of interest (positively, negatively,
826 or without restriction) were averaged over the four 250 ms bins corresponding to the Late Trace
827 period, and then we subtracted the shuffled from the unshuffled fraction to account for odor
828 coding.

829 *Changes relative to Baseline*

830 In order to assess changes in neural activity relative to the Baseline period, we first grouped all
831 Unrewarded (Nothing) and Rewarded (Fixed and Variable) trials for each neuron. We then ran a
832 rank-sum test between Late Trace activity and Baseline activity, separately on each neuron and
833 trial type grouping. Finally, we computed the fraction of cells per mouse that increased or
834 decreased significantly ($\alpha = 0.05$) and then ran paired t -tests on the respective fractions for
835 Rewarded versus Unrewarded trials types.

836 *Comparisons across subregions, hemispheres, and genotypes*

837 Whenever subregions, hemispheres, or genotypes were directly compared, we randomly
838 subsampled the number of neurons so that population sizes were identical across this
839 comparison. For subregion and hemisphere (lesioned vs. control), this matching was done
840 within-animal. When comparing subregions, we excluded a subregion from an animal if it did
841 not contain at least 40 neurons, hence the differing number of dots (animals) per subregion
842 (Extended Data Fig. 2e, 3e, 6a-d, 6f. For genotype (D1 vs. D2 MSNs), matching was done
843 across-animals, for the entire population of D1 or D2 neurons. To allow for higher neuron
844 counts, all of these decoding analyses were performed on pseudo-populations.

845 *Artificial neural network-based distribution decoding*

846 To determine whether neural populations contained sufficient information to reconstruct the
847 complete reward distribution, rather than simply perform binary classification based on reward
848 variance, we constructed an artificial neural network (ANN)-based distribution decoder. Pseudo-

849 population activity from the distribution-coding subpopulation r was first mapped into 16
850 dimensions by a trainable, unregularized decoding matrix W . The network takes Wr as input and
851 outputs the predicted distribution. It has one input layer, two hidden layers, and one output layer.
852 Each of the two hidden layers had 32 neurons and used the non-linear activation function $f(x) =$
853 $\ln(1 + \exp(x + I)) - I$, which is close to the identity function for $x \gg 0$ and to -1 for $x \ll 0$.
854 The output layer had size 4, with each dimension corresponding to a possible reward size (0, 2,
855 4, or 6 μL). After linear combination, we also applied the nonlinear function $f(x)$ as specified
856 above, followed by the softmax function to turn the output into a normalized probability
857 distribution.

858 We applied stochastic gradient descent (SGD) to minimize the following loss function based on
859 the 1-Wasserstein distance (D):

$$860 \quad L(W, \text{network weights}) = \langle D(\text{decoded_dist}, \text{groundtruth_dist}) \rangle + \lambda \|\text{network weights}\|_2^2,$$

861 where D is defined as $D(P, Q) = \sum_n |P(r_n) - Q(r_n)|$ for discrete cumulative distribution
862 functions (CDFs) P and Q , where the sum is over all used reward magnitudes, and where r_n is the
863 respective reward magnitude. In other words, the 1-Wasserstein distance measures the unsigned
864 area between two CDFs. For plotting, we normalized this metric by dividing by the minimum
865 achievable Wasserstein distance that would result from predicting the same distribution for every
866 trial type across the training and test sets (“Wasserstein distance relative to reference”).

867 For all experiments, λ was set to 0.02 and the learning rate was 0.002. All the trainable weights
868 were randomly initialized with a mean of 0 and standard deviation of 1, and then divided by 15.
869 For each disjoint pseudo-population, we trained each of 5 candidate ANNs initialized randomly
870 and differently for 1,200 iterations, and picked the best-performing one to further train for
871 10,000 iterations. The ANN was implemented in Julia (v. 1.6.7) and trained on a GPU (NVIDIA,
872 GeForce RTX 2070).

873 In the standard decoding setting, all six trial types were included in the training and testing sets
874 (with different trials in each). For decoding restricted to trial types with the same mean, only
875 Fixed and Variable trial types were used, but split according to the same logic. In both cases, we
876 performed decoding independently from each mouse, and we compared our results to what
877 happened when we randomly shuffled the odor-distribution mappings before training. If merely
878 odor identity (or, in the restricted case, mean) is encoded, then the ordered and shuffled networks
879 should attain similar performance.

880 Finally, in the transfer analysis, in a similar spirit to CCGP, we trained on only four trial types
881 and then tested on the held-out two trial types. “Matched” transfers used one Fixed and one
882 Variable odor in the training set, assigned to the proper distribution, and evaluated performance
883 on the corresponding test odor. “Mismatched” transfers used either two Fixed or two Variable
884 odors in the training set, assigning one to each distribution, and evaluated performance on the

885 held-out odors, again assigning one to each distribution. Nothing trial types were always
886 assigned to Nothing distributions. To gain statistical power, we pooled neurons across mice for
887 these analyses.

888 Computational Modeling

889 In this section, we briefly review the theory behind various distributional RL algorithms before
890 specifying the details of our implementation, for the purpose of comparing the learned code to
891 neural activity and generating predictions for optogenetic perturbations. All models were trained
892 for 2,000 trials per distribution.

893 *Reflected expectile distributional RL (REDRL)*

894 EDRL was first put forward as a novel machine learning algorithm⁷⁶ and later used to explain
895 dopamine neuron diversity in the mammalian midbrain⁴. EDRL approximately minimizes the
896 expectile regression loss function (ER):

$$897 \quad ER(V; \mu, \tau) = \mathbb{E}_{Z \sim \mu} [\tau \mathbb{1}_{Z > V} + (1 - \tau) \mathbb{1}_{Z \leq V}] (Z - V)^2,$$

898 where V is the value predictor, μ is the target distribution, Z is a random sample from μ , τ is the
899 asymmetry, and $\mathbb{1}$ is the indicator function, which is 1 when the subscript is satisfied and 0 when
900 it is violated. It is an asymmetrically-weighted squared error loss function; in this sense, it
901 generalizes the mean (squared error loss, equivalent to the 0.5th expectile) just as quantiles
902 generalize the median⁹⁷.

903 EDRL and REDRL minimize this ER loss function simultaneously for many values of τ , indexed
904 by i , generally using SGD with respect to the value predictors (or their parameters). This
905 formulation is sufficiently general that it can be combined with nonlinear function approximation
906 and temporal difference learning methods, and its effectiveness has been demonstrated on the
907 suite of Atari video games⁷⁶. However, for simplicity, here we present the Rescorla-Wagner¹⁷³
908 version of the update rule for tabular states, so the random sample from μ reduces to simply the
909 reward, r . This is the learning rule depicted in Fig. 3:

$$910 \quad \delta_i = r - V_i$$

$$911 \quad V_i \leftarrow V_i + \alpha_i^- \cdot \delta_i, \text{ if } \delta_i \leq 0$$

$$912 \quad V_i \leftarrow V_i + \alpha_i^+ \cdot \delta_i, \text{ if } \delta_i > 0$$

913 For the learning simulations (Fig. 3a), we used $\alpha = \alpha_i^+ + \alpha_i^- = 0.03$ and initialized all value
914 predictors to 2.

915 In the biological implementation of the REDRL algorithm (Fig. 3d-g), we decompose this update

916 into two piecewise linear functions. The first function models dopamine RPEs, which are
917 allowed to take on different slopes in the positive and negative domains, α_i^+ and α_i^- . The second
918 function differs between D1 and D2 MSNs (indexed by m) by a reflection over the y-axis. It
919 maps changes in dopamine firing into changes in synaptic weights³⁹, which we'll parameterize
920 here by $\beta_m^{-/+}$ (equal to 0.75/3 for D1 and 3/0.75 for D2 MSNs for the purpose of Fig. 3).

921 Composing these functions gives rise to the following update rules:

922
$$D1_i \leftarrow D1_i + \alpha_i'^- \cdot \beta_{D1}^- \cdot \delta_i, \text{ if } \delta_i \leq 0$$

923
$$D1_i \leftarrow D1_i + \alpha_i'^+ \cdot \beta_{D1}^+ \cdot \delta_i, \text{ if } \delta_i > 0$$

924
$$D2_i \leftarrow D2_i - \alpha_i'^- \cdot \beta_{D2}^- \cdot \delta_i, \text{ if } \delta_i \leq 0$$

925
$$D2_i \leftarrow D2_i - \alpha_i'^+ \cdot \beta_{D2}^+ \cdot \delta_i, \text{ if } \delta_i > 0$$

926 Note that D1 and D2 neurons receive unique indices i , so there is no overlap in the idealized
927 case. As a consequence of the opponent plasticity rule, changes in synaptic weights in D1 and D2
928 MSNs have opposing effects on the encoded value predictor, modeled simply by the identity
929 function (for D1 MSNs) or its negation, (for D2 MSNs). Therefore, this update rule becomes
930 equivalent to the algorithmic rule if we let $\alpha_i^- = \alpha_i'^- \cdot \beta_m^-$ and $\alpha_i^+ = \alpha_i'^+ \cdot \beta_m^+$. The degree of
931 optimism or pessimism is parameterized by the dimensionless quantity $\tau_i = \frac{\alpha_i^+}{\alpha_i^+ + \alpha_i^-}$, which
932 ranges from 0 to 1. Importantly, τ_i uses the net asymmetries learned by the MSNs as opposed to
933 the asymmetries of the dopamine neurons. Therefore, both the expectile that is learned in the
934 striatum and the zero-crossing point of the corresponding dopamine neuron are dictated by τ_i ,
935 which can give rise to multiple dopamine neurons with the same apparent asymmetry but
936 different zero-crossing points. This stands in contrast to the EDRL model, in which the dopamine
937 neuron asymmetries alone fully determine the zero-crossing point, but nonetheless predicts the
938 observed correlation between zero-crossing points and asymmetries⁴.

939 For D1 MSNs $\beta_m^+ > \beta_m^-$ and so τ_i skews optimistic; analogously, for D2 MSNs $\beta_m^+ < \beta_m^-$, and τ_i
940 skews pessimistic. The precise distribution of τ 's will depend on the distribution of dopamine
941 neuron asymmetries (α_i^+ and α_i^-) as well as the ratio of β_m^+ to β_m^- , neither of which has been
942 measured precisely. To avoid making too many assumptions and to simplify interpretation, we
943 plotted all REDRL results based on a simulation of 10 predictors with uniform spacing of τ_i
944 between 0.05 and 0.95, with all $\tau_i > 0.5$ assigned to D1 MSNs and all $\tau_i < 0.5$ assigned to D2
945 MSNs. Furthermore, we directly computed the expectiles of the relevant reward distributions
946 (rather than obtaining them incrementally from samples and updates) in order to eliminate noise.
947 We confirmed that all of our main results were robust to these choices of τ and simulation
948 approach.

949 *Quantile distributional RL (QDRL)*

950 QDRL is exactly akin to EDRL, except that we minimize the quantile regression (QR) loss⁷²:

951
$$QR(V; \mu, \tau) = \mathbb{E}_{Z \sim \mu} [[\tau \mathbb{1}_{Z > V} + (1 - \tau) \mathbb{1}_{Z < V}] |Z - V|],$$

952 This is an asymmetrically-weighted absolute value loss function, which would return the median
953 when positive and negative errors are balanced ($\tau = 0.5$). The update rule, derived by SGD,
954 utilizes only the sign of the prediction error, not its magnitude⁹⁷:

955
$$V_i \leftarrow V_i - \alpha_i^-, \text{ if } \delta_i < 0$$

956
$$V_i \leftarrow V_i + \alpha_i^+, \text{ if } \delta_i > 0$$

957 Unlike expectiles, quantiles have an intuitive interpretation: the τ -th quantile is the number such
958 that τ fraction of samples from the distribution fall below that value and $1 - \tau$ fall above it. It is
959 therefore the inverse of the cumulative distribution function (CDF). We additionally
960 implemented a “reflected” version of QDRL by applying the same transformation to D2 MSNs,
961 those predictors with $\tau_i < 0.5$.

962 We also note that it is possible to interpolate between EDRL and QDRL using Huber
963 quantiles^{72,174}. This is simply an asymmetric squared loss within a certain interval (controlled by
964 a hyperparameter κ), and a standard quantile loss outside this interval. The update rule is likewise
965 a combination of EDRL and QDRL: piecewise linear within some range before saturating. This
966 rule would obtain if, for example, plasticity could only change some maximum amount in either
967 direction at any given time, as is likely the case in the brain. Notably, the Huber quantile loss is
968 frequently used in machine learning applications⁷².

969 *Categorical distributional RL (CDRL)*

970 CDRL⁷¹ adopts a very different approach to learning the reward distribution. Rather than a
971 quantile or expectile function, CDRL imagines a set of “atoms”, which function similarly to bins
972 of a histogram. For that reason, we model these “categorical codes” using one hypothetical
973 neuron per reward size (0–8 μL), in increments of 2 μL . The height of that bin is then assumed to
974 be linearly (and positively) related to the firing rate of that neuron. Generalizing this scheme to
975 use basis functions over bin values does not qualitatively alter the predictions.

976 *Laplace and cumulative code*

977 The Laplace code⁸³ grew out of an effort to devise a fully local temporal difference (TD)
978 learning rule for distributional RL. Its teaching signal is simply a sigmoidal function of reward: if
979 reward exceeds some threshold, the neuron fires, and thresholds are heterogeneous across
980 neurons. In the limit of infinitely steep sigmoids (Heaviside step functions), the value predictors

981 converge to the probability that the reward exceeds the given threshold (discounted and summed
982 over future time steps, in the TD case). This exceedance probability is equal to $1 - \text{CDF}$ of the
983 reward distribution, for our simplified Rescorla-Wagner setting. By analogy to CDRL, we chose
984 to model neural activity as linearly and positively related to this value of $1 - \text{CDF}$ at each of the
985 reward bins. For completeness, we also investigated a “cumulative” code, which was just the
986 CDF at each reward bin, or $1 -$ the Laplace code. The spatial derivative of this cumulative code
987 is then equivalent to the categorical code, assuming sufficient support.

988 *Actor Uncertainty (AU) model*

989 The AU model⁶⁶ manages to learn about reward uncertainty using biologically-plausible learning
990 rules in D1 and D2 MSNs. We therefore wanted to test its predictions against these other models.
991 The AU model makes use of two value predictors: one D1 and one D2 MSN, which learn as
992 follows:

$$993 \quad V = D1 - D2$$

$$994 \quad D1 \leftarrow D1 + \alpha |r - V|_+ - \beta \cdot D1$$

$$995 \quad D2 \leftarrow D2 + \alpha |r - V|_- - \beta \cdot D2$$

996 Here, $|x|_+ = \max(x, 0)$ and $|x|_- = \max(-x, 0)$, and $0 < \beta < 1$ scales the decay term to ensure
997 stability. Using this model, it can be shown⁶⁶ that $D1 - D2$ encodes an estimate of mean reward,
998 and $D1 + D2$ encodes an estimate of reward spread. For our implementation, we set $\alpha = 0.1$ and β
999 $= 0.01$.

1000 *Distributed AU model*

1001 The distributed AU model¹⁷⁵ works similarly, except that we now allow there to be different
1002 learning rates α_i^+ and α_i^- for D1 and D2 MSNs, respectively, just as in the distributional RL
1003 setting. The difference $V_i = D1_i - D2_i$ approximates the τ_i -th expectile, biased by β . For our
1004 simulations, we chose $\alpha = \alpha_i^+ + \alpha_i^- = 0.2$ and $\beta = 0.01$.

1005 *Modeling perturbations*

1006 Simulating optogenetic inhibition and excitation in these models (Extended Data Fig. 11)
1007 required slightly different choices, depending on the type of code. For expectile, quantile, and
1008 AU-based models, we clamped the relevant simulated neuron(s) to either 0 or 8, the maximum
1009 reward value across all distributions, to simulate model inhibition and excitation, respectively.
1010 Note that it was the neural activity ($D1_i$ or $D2_i$) that we were directly clamping when applicable,
1011 not the value prediction it encoded (V_i). For the expectile and quantile models, optimistic and
1012 pessimistic perturbations meant clamping the value of predictors with $\tau_i > 0.5$ and $\tau_i < 0.5$
1013 respectively. For the AU model, they were identified with the D1 and D2 MSN, respectively.

1014 Finally, for the distributed AU model, we implemented two versions of the perturbation, one in
1015 which all D1 (optimistic) or all D2 (pessimistic) neurons were manipulated, and one in which
1016 only those with $\tau_i > 0.5$ or $\tau_i < 0.5$, respectively, were manipulated. We call the latter the “Partial
1017 Distributed AU” model, for the purposes of model comparison. For the AU models, it is only the
1018 difference $D1_i - D2_i$ that is bounded within the range of reward sizes, not the activities
1019 individually. We therefore added or subtracted a fixed amount (the maximum reward size across
1020 all trial types, 8 μL) across reward predictors to simulate excitation or inhibition, respectively, in
1021 these models, rather than clamping their value to a constant.

1022 For categorical, cumulative, and Laplace codes, the semantics of each simulated neuron are
1023 different: their activations range from 0 to 1 and encode a (cumulative) probability, rather than a
1024 value. Thus, inhibiting or exciting them meant changing the relevant probability to 0 or 1,
1025 respectively. Pessimistic neurons were those that corresponded to the 0 or 2 μL bins, and
1026 optimistic neurons corresponded to 6 and 8 μL . To reconstitute a properly-normalized probability
1027 distribution after the perturbation, in the case of the categorical code, we divided by the sum of
1028 the predictors (or made it a uniform distribution if the sum was zero). For the categorical and
1029 Laplace codes, we took the spatial derivative of the implied CDF, subtracted off the minimum if
1030 any value was negative, and then divided by the sum (or made it uniform if the sum was zero).

1031 In all cases, we found the mean of the (imputed) perturbed probability distribution and then
1032 compared it to the mean without any perturbation to model the effect of optogenetic
1033 manipulation on lick rate.

1034 *Model comparison*

1035 We used the predicted Manipulation – No Manipulation differences from each model as a
1036 regressor with which to predict the difference in licking across trial types, averaged across mice,
1037 using linear regression (with no intercept term). Separate regressions were fit for inhibition and
1038 excitation to allow for potentially different scaling in each case, and their coefficients of
1039 determination were averaged to produce a single summary measure of goodness of fit.

1040 *Data availability*

1041 Pre-processed data will be posted to online repositories upon publication.

1042 *Code availability*

1043 Analysis code will be posted to online repositories upon publication.

1044 **References**

- 1045 1. Bellemare, M. G., Dabney, W. & Rowland, M. *Distributional Reinforcement Learning*.
1046 (MIT Press, 2023).
- 1047 2. Schultz, W., Dayan, P. & Montague, P. R. A neural substrate of prediction and reward.
1048 *Science* **275**, 1593–1599 (1997).
- 1049 3. Doya, K. Reinforcement learning: Computational theory and biological mechanisms. *HFSP*
1050 *J.* **1**, 30–40 (2007).
- 1051 4. Dabney, W. *et al.* A distributional code for value in dopamine-based reinforcement learning.
1052 *Nature* **577**, 671–675 (2020).
- 1053 5. Shin, J. H., Kim, D. & Jung, M. W. Differential coding of reward and movement
1054 information in the dorsomedial striatal direct and indirect pathways. *Nat. Commun.* **9**, 404
1055 (2018).
- 1056 6. Nonomura, S. *et al.* Monitoring and Updating of Action Selection for Goal-Directed
1057 Behavior through the Striatal Direct and Indirect Pathways. *Neuron* **99**, 1302-1314.e5
1058 (2018).
- 1059 7. Hikida, T., Kimura, K., Wada, N., Funabiki, K. & Nakanishi, S. Distinct Roles of Synaptic
1060 Transmission in Direct and Indirect Striatal Pathways to Reward and Aversive Behavior.
1061 *Neuron* **66**, 896–907 (2010).
- 1062 8. Kravitz, A. V. *et al.* Regulation of parkinsonian motor behaviours by optogenetic control of
1063 basal ganglia circuitry. *Nature* **466**, 622–626 (2010).
- 1064 9. Lobo, M. K. *et al.* Cell type-specific loss of BDNF signaling mimics optogenetic control of
1065 cocaine reward. *Science* **330**, 385–390 (2010).
- 1066 10. Kravitz, A. V., Tye, L. D. & Kreitzer, A. C. Distinct roles for direct and indirect pathway
1067 striatal neurons in reinforcement. *Nat. Neurosci.* **15**, 816–818 (2012).
- 1068 11. Tai, L.-H., Lee, A. M., Benavidez, N., Bonci, A. & Wilbrecht, L. Transient stimulation of
1069 distinct subpopulations of striatal neurons mimics changes in action value. *Nat. Neurosci.*
1070 **15**, 1281–1289 (2012).
- 1071 12. Hikida, T. *et al.* Pathway-specific modulation of nucleus accumbens in reward and aversive
1072 behavior via selective transmitter receptors. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 342–347
1073 (2013).
- 1074 13. Yttri, E. A. & Dudman, J. T. Opponent and bidirectional control of movement velocity in
1075 the basal ganglia. *Nature* **533**, 402–406 (2016).
- 1076 14. Parker, J. G. *et al.* Diametric neural ensemble dynamics in parkinsonian and dyskinetic
1077 states. *Nature* **557**, 177–182 (2018).
- 1078 15. Cruz, B. F. *et al.* Action suppression reveals opponent parallel control via striatal circuits.
1079 *Nature* **607**, 521–526 (2022).
- 1080 16. Grillner, S., Robertson, B. & Stephenson-Jones, M. The evolutionary origin of the vertebrate
1081 basal ganglia and its role in action selection. *J. Physiol.* **591**, 5425–5431 (2013).
- 1082 17. Floresco, S. B. The nucleus accumbens: an interface between cognition, emotion, and
1083 action. *Annu. Rev. Psychol.* **66**, 25–52 (2015).

- 1084 18. Cox, J. & Witten, I. B. Striatal circuits for reward learning and decision-making. *Nat. Rev.*
1085 *Neurosci.* **20**, 482–494 (2019).
- 1086 19. Montague, P. R., Dayan, P. & Sejnowski, T. J. A Framework for Mesencephalic Dopamine
1087 Systems Based on Predictive Hebbian Learning. *The Journal of Neuroscience* **16**, 1936–
1088 1947 (1996).
- 1089 20. Cohen, J. Y., Haesler, S., Vong, L., Lowell, B. B. & Uchida, N. Neuron-type-specific
1090 signals for reward and punishment in the ventral tegmental area. *Nature* **482**, 85–88 (2012).
- 1091 21. Sutton, R. S. & Barto, A. G. *Reinforcement Learning: An Introduction*. vol. 2 (MIT Press,
1092 2018).
- 1093 22. Yagishita, S. *et al.* A critical time window for dopamine actions on the structural plasticity
1094 of dendritic spines. *Science* **345**, 1616–1620 (2014).
- 1095 23. Schultz, W., Apicella, P., Scarnati, E. & Ljungberg, T. Neuronal activity in monkey ventral
1096 striatum related to the expectation of reward. *J. Neurosci.* **12**, 4595–4610 (1992).
- 1097 24. Taha, S. A. & Fields, H. L. Encoding of palatability and appetitive behaviors by distinct
1098 neuronal populations in the nucleus accumbens. *J. Neurosci.* **25**, 1193–1202 (2005).
- 1099 25. Roesch, M. R., Singh, T., Brown, P. L., Mullins, S. E. & Schoenbaum, G. Ventral striatal
1100 neurons encode the value of the chosen action in rats deciding between differently delayed
1101 or sized rewards. *J. Neurosci.* **29**, 13365–13376 (2009).
- 1102 26. Ito, M. & Doya, K. Distinct neural representation in the dorsolateral, dorsomedial, and
1103 ventral parts of the striatum during fixed- and free-choice tasks. *J. Neurosci.* **35**, 3499–3514
1104 (2015).
- 1105 27. Strait, C. E., Slezzer, B. J. & Hayden, B. Y. Signatures of Value Comparison in Ventral
1106 Striatum Neurons. *PLoS Biol.* **13**, e1002173 (2015).
- 1107 28. Shin, E. J. *et al.* Robust and distributed neural representation of action values. *eLife* **10**,
1108 (2021).
- 1109 29. Ottenheimer, D., Richard, J. M. & Janak, P. H. Ventral pallidum encodes relative reward
1110 value earlier and more robustly than nucleus accumbens. *Nat. Commun.* **9**, 4350 (2018).
- 1111 30. Lee, D., Liu, L. & Root, C. M. Transformation of value signaling in a striatopallidal circuit.
1112 *eLife* (2023) doi:10.7554/elife.90976.
- 1113 31. Padoa-Schioppa, C. & Assad, J. A. Neurons in the orbitofrontal cortex encode economic
1114 value. *Nature* **441**, 223–226 (2006).
- 1115 32. Hattori, R., Danskin, B., Babic, Z., Mlynaryk, N. & Komiyama, T. Area-Specificity and
1116 Plasticity of History-Dependent Value Coding During Learning. *Cell* **177**, 1858–1872.e15
1117 (2019).
- 1118 33. Enel, P., Wallis, J. D. & Rich, E. L. Stable and dynamic representations of value in the
1119 prefrontal cortex. *eLife* **9**, (2020).
- 1120 34. Ottenheimer, D. J., Hjort, M. M., Bowen, A. J., Steinmetz, N. A. & Stuber, G. D. A stable,
1121 distributed code for cue value in mouse cortex during reward learning. *eLife* **12**, (2023).
- 1122 35. Bao, C. *et al.* The rat frontal orienting field dynamically encodes value for economic
1123 decisions under risk. *Nat. Neurosci.* **26**, 1942–1952 (2023).

- 1124 36. Bari, B. A. *et al.* Stable Representations of Decision Variables for Flexible Behavior.
1125 *Neuron* **103**, 922–933.e7 (2019).
- 1126 37. Hirokawa, J., Vaughan, A., Masset, P., Ott, T. & Kepecs, A. Frontal cortex neuron types
1127 categorically encode single decision variables. *Nature* **576**, 446–451 (2019).
- 1128 38. Gerfen, C. R. & Surmeier, D. J. Modulation of striatal projection systems by dopamine.
1129 *Annu. Rev. Neurosci.* **34**, 441–466 (2011).
- 1130 39. Iino, Y. *et al.* Dopamine D2 receptors in discrimination learning and spine enlargement.
1131 *Nature* **579**, 555–560 (2020).
- 1132 40. Lee, S. J. *et al.* Cell-type-specific asynchronous modulation of PKA by dopamine in
1133 learning. *Nature* **590**, 451–456 (2021).
- 1134 41. Cui, G. *et al.* Concurrent activation of striatal direct and indirect pathways during action
1135 initiation. *Nature* **494**, 238–242 (2013).
- 1136 42. Tecuapetla, F., Matias, S., Dugue, G. P., Mainen, Z. F. & Costa, R. M. Balanced activity in
1137 basal ganglia projection pathways is critical for contraversive movements. *Nat. Commun.* **5**,
1138 4315 (2014).
- 1139 43. Klaus, A. *et al.* The Spatiotemporal Organization of the Striatum Encodes Action Space.
1140 *Neuron* **95**, 1171–1180.e7 (2017).
- 1141 44. Markowitz, J. E. *et al.* The Striatum Organizes 3D Behavior via Moment-to-Moment Action
1142 Selection. *Cell* **174**, 44–58 (2018).
- 1143 45. Tan, B. *et al.* Dynamic processing of hunger and thirst by common mesolimbic neural
1144 ensembles. *Proc. Natl. Acad. Sci. U. S. A.* **119**, e2211688119 (2022).
- 1145 46. Menegas, W., Babayan, B. M., Uchida, N. & Watabe-Uchida, M. Opposite initialization to
1146 novel cues in dopamine signaling in ventral and posterior striatum in mice. *eLife* **6**, (2017).
- 1147 47. Menegas, W., Akiti, K., Amo, R., Uchida, N. & Watabe-Uchida, M. Dopamine neurons
1148 projecting to the posterior striatum reinforce avoidance of threatening stimuli. *Nat.*
1149 *Neurosci.* **21**, 1421–1430 (2018).
- 1150 48. Watabe-Uchida, M. & Uchida, N. Multiple Dopamine Systems: Weal and Woe of
1151 Dopamine. *Cold Spring Harb. Symp. Quant. Biol.* **83**, 83–95 (2018).
- 1152 49. de Jong, J. W. *et al.* A Neural Circuit Mechanism for Encoding Aversive Stimuli in the
1153 Mesolimbic Dopamine System. *Neuron* **101**, 133–151.e7 (2019).
- 1154 50. Tsutsui-Kimura, I. *et al.* Distinct temporal difference error signals in dopamine axons in
1155 three regions of the striatum in a decision-making task. *eLife* **9**, (2020).
- 1156 51. Engelhard, B. *et al.* Specialized coding of sensory, motor and cognitive variables in VTA
1157 dopamine neurons. *Nature* **570**, 509–513 (2019).
- 1158 52. Akiti, K. *et al.* Striatal dopamine explains novelty-induced behavioral dynamics and
1159 individual variability in threat prediction. *Neuron* **110**, 3789–3804.e9 (2022).
- 1160 53. Lee, R. S., Sagiv, Y., Engelhard, B., Witten, I. B. & Daw, N. D. A feature-specific
1161 prediction error model explains dopaminergic heterogeneity. *bioRxiv* 2022.02.28.482379
1162 (2023) doi:10.1101/2022.02.28.482379.
- 1163 54. Jeong, H. *et al.* Mesolimbic dopamine release conveys causal associations. *Science* **378**,

- 1164 eabq6740 (2022).
- 1165 55. Coddington, L. T., Lindo, S. E. & Dudman, J. T. Mesolimbic dopamine adapts the rate of
1166 learning from action. *Nature* **614**, 294–302 (2023).
- 1167 56. Costa, V. D., Dal Monte, O., Lucas, D. R., Murray, E. A. & Averbeck, B. B. Amygdala and
1168 Ventral Striatum Make Distinct Contributions to Reinforcement Learning. *Neuron* **92**, 505–
1169 517 (2016).
- 1170 57. Rothenhoefer, K. M. *et al.* Effects of Ventral Striatum Lesions on Stimulus-Based versus
1171 Action-Based Reinforcement Learning. *J. Neurosci.* **37**, 6902–6914 (2017).
- 1172 58. Floresco, S. B., Montes, D. R., Tse, M. M. T. & van Holstein, M. Differential Contributions
1173 of Nucleus Accumbens Subregions to Cue-Guided Risk/Reward Decision Making and
1174 Implementation of Conditional Rules. *J. Neurosci.* **38**, 1901–1914 (2018).
- 1175 59. St Onge, J. R. & Floresco, S. B. Dopaminergic modulation of risk-based decision making.
1176 *Neuropsychopharmacology* **34**, 681–697 (2009).
- 1177 60. St Onge, J. R., Chiu, Y. C. & Floresco, S. B. Differential effects of dopaminergic
1178 manipulations on risky choice. *Psychopharmacology* **211**, 209–221 (2010).
- 1179 61. Zalocusky, K. A. *et al.* Nucleus accumbens D2R cells signal prior outcomes and control
1180 risky decision-making. *Nature* **531**, 642–646 (2016).
- 1181 62. Mortazavi, L. *et al.* D2/3 Agonist during Learning Potentiates Cued Risky Choice. *J.*
1182 *Neurosci.* **43**, 979–992 (2023).
- 1183 63. Yager, L. M., Garcia, A. F., Wunsch, A. M. & Ferguson, S. M. The ins and outs of the
1184 striatum: role in drug addiction. *Neuroscience* **301**, 529–541 (2015).
- 1185 64. Everitt, B. J. & Robbins, T. W. Drug Addiction: Updating Actions to Habits to Compulsions
1186 Ten Years On. *Annu. Rev. Psychol.* **67**, 23–50 (2016).
- 1187 65. Gatto, E. M. & Aldinio, V. Impulse Control Disorders in Parkinson’s Disease. A Brief and
1188 Comprehensive Review. *Front. Neurol.* **10**, 351 (2019).
- 1189 66. Mikhael, J. G. & Bogacz, R. Learning Reward Uncertainty in the Basal Ganglia. *PLoS*
1190 *Comput. Biol.* **12**, e1005062 (2016).
- 1191 67. Ma, W. J., Beck, J. M., Latham, P. E. & Pouget, A. Bayesian inference with probabilistic
1192 population codes. *Nat. Neurosci.* **9**, 1432–1438 (2006).
- 1193 68. Vértes, E. & Sahani, M. Flexible and accurate inference and learning for deep generative
1194 models. in *Advances in Neural Information Processing Systems* (eds. Bengio, S. et al.) vol.
1195 31 4166–4175 (Curran Associates, Inc., 2018).
- 1196 69. Walker, E. Y., Cotton, R. J., Ma, W. J. & Tolias, A. S. A neural basis of probabilistic
1197 computation in visual cortex. *Nat. Neurosci.* **23**, 122–129 (2020).
- 1198 70. Walker, E. Y. *et al.* Studying the neural representations of uncertainty. *Nat. Neurosci.* **26**,
1199 1857–1867 (2023).
- 1200 71. Bellemare, M. G., Dabney, W. & Munos, R. A Distributional Perspective on Reinforcement
1201 Learning. in *Proceedings of the 34th International Conference on Machine Learning* (eds.
1202 Precup, D. & Teh, Y. W.) vol. 70 449–458 (PMLR, 06–11 Aug 2017).
- 1203 72. Dabney, W., Rowland, M., Bellemare, M. & Munos, R. Distributional Reinforcement

- 1204 Learning With Quantile Regression. in *Proceedings of the AAAI Conference on Artificial*
1205 *Intelligence* vol. 32 (2018).
- 1206 73. Wurman, P. R. *et al.* Outracing champion Gran Turismo drivers with deep reinforcement
1207 learning. *Nature* **602**, 223–228 (2022).
- 1208 74. Eshel, N., Tian, J., Bukwich, M. & Uchida, N. Dopamine neurons share common response
1209 function for reward prediction error. *Nat. Neurosci.* **19**, 479–486 (2016).
- 1210 75. Rothenhoefer, K. M., Hong, T., Alikaya, A. & Stauffer, W. R. Rare rewards amplify
1211 dopamine responses. *Nat. Neurosci.* **24**, 465–469 (2021).
- 1212 76. Rowland, M. *et al.* Statistics and Samples in Distributional Reinforcement Learning. in
1213 *Proceedings of the 36th International Conference on Machine Learning* (eds. Chaudhuri, K.
1214 & Salakhutdinov, R.) vol. 97 5528–5536 (PMLR, 09–15 Jun 2019).
- 1215 77. Lyle, C., Castro, P. S. & Bellemare, M. G. A Comparative Analysis of Expected and
1216 Distributional Reinforcement Learning. *arXiv [cs.LG]* (2019).
- 1217 78. Such, F. P. *et al.* An Atari model zoo for analyzing, visualizing, and comparing deep
1218 reinforcement learning agents. in *Proceedings of the Twenty-Eighth International Joint*
1219 *Conference on Artificial Intelligence* 3260–3267 (International Joint Conferences on
1220 Artificial Intelligence Organization, 2019).
- 1221 79. Nikolov, N., Kirschner, J., Berkenkamp, F. & Krause, A. Information-Directed Exploration
1222 for Deep Reinforcement Learning. *arXiv [cs.LG]* (2018).
- 1223 80. Mavrin, B. *et al.* Distributional Reinforcement Learning for Efficient Exploration. *arXiv*
1224 *[cs.LG]* (2019).
- 1225 81. Clements, W. R., Van Delft, B., Robaglia, B.-M., Slaoui, R. B. & Toth, S. Estimating Risk
1226 and Uncertainty in Deep Reinforcement Learning. *arXiv [cs.LG]* (2019).
- 1227 82. Zhang, S. & Yao, H. QUOTA: The Quantile Option Architecture for Reinforcement
1228 Learning. *AAAI* **33**, 5797–5804 (2019).
- 1229 83. Tano, P., Dayan, P. & Pouget, A. A local temporal difference code for distributional
1230 reinforcement learning. in *Advances in Neural Information Processing Systems* (eds.
1231 Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F. & Lin, H.) vol. 33 13662–13673
1232 (2020).
- 1233 84. Louie, K. Asymmetric and adaptive reward coding via normalized reinforcement learning.
1234 *PLoS Comput. Biol.* **18**, e1010350 (2022).
- 1235 85. Schütt, H. H., Kim, D. & Ma, W. J. Reward prediction error neurons implement an efficient
1236 code for reward. *bioRxiv* 2022.11.03.515104 (2023) doi:10.1101/2022.11.03.515104.
- 1237 86. McCoy, A. N. & Platt, M. L. Risk-sensitive neurons in macaque posterior cingulate cortex.
1238 *Nat. Neurosci.* **8**, 1220–1227 (2005).
- 1239 87. O’Neill, M. & Schultz, W. Coding of reward risk by orbitofrontal neurons is mostly distinct
1240 from coding of reward value. *Neuron* **68**, 789–800 (2010).
- 1241 88. Monosov, I. E. Anterior cingulate is a source of valence-specific information about value
1242 and uncertainty. *Nat. Commun.* **8**, 134 (2017).
- 1243 89. Monosov, I. E. & Hikosaka, O. Selective and graded coding of reward uncertainty by

- 1244 neurons in the primate anterodorsal septal region. *Nat. Neurosci.* **16**, 756–762 (2013).
- 1245 90. Monosov, I. E., Leopold, D. A. & Hikosaka, O. Neurons in the Primate Medial Basal
1246 Forebrain Signal Combined Information about Reward Uncertainty, Value, and Punishment
1247 Anticipation. *J. Neurosci.* **35**, 7443–7459 (2015).
- 1248 91. Ledbetter, N. M., Chen, C. D. & Monosov, I. E. Multiple Mechanisms for Processing
1249 Reward Uncertainty in the Primate Basal Forebrain. *J. Neurosci.* **36**, 7852–7864 (2016).
- 1250 92. Stringer, C. *et al.* Spontaneous behaviors drive multidimensional, brainwide activity.
1251 *Science* **364**, 255 (2019).
- 1252 93. Tian, J. *et al.* Distributed and Mixed Information in Monosynaptic Inputs to Dopamine
1253 Neurons. *Neuron* **91**, 1374–1389 (2016).
- 1254 94. Orbán, G., Berkes, P., Fiser, J. & Lengyel, M. Neural Variability and Sampling-Based
1255 Probabilistic Representations in the Visual Cortex. *Neuron* **92**, 530–543 (2016).
- 1256 95. Echeveste, R., Aitchison, L., Hennequin, G. & Lengyel, M. Cortical-like dynamics in
1257 recurrent circuits optimized for sampling-based probabilistic inference. *Nat. Neurosci.* **23**,
1258 1138–1149 (2020).
- 1259 96. Bernardi, S. *et al.* The Geometry of Abstraction in the Hippocampus and Prefrontal Cortex.
1260 *Cell* **183**, 954–967.e21 (2020).
- 1261 97. Lowet, A. S., Zheng, Q., Matias, S., Drugowitsch, J. & Uchida, N. Distributional
1262 Reinforcement Learning in the Brain. *Trends Neurosci.* **43**, 980–997 (2020).
- 1263 98. Rice, M. E. & Cragg, S. J. Dopamine spillover after quantal release: rethinking dopamine
1264 transmission in the nigrostriatal pathway. *Brain Res. Rev.* **58**, 303–313 (2008).
- 1265 99. Dreyer, J. K., Herrik, K. F., Berg, R. W. & Hounsgaard, J. D. Influence of phasic and tonic
1266 dopamine release on receptor activation. *J. Neurosci.* **30**, 14273–14283 (2010).
- 1267 100. Faust, T. W., Mohebi, A. & Berke, J. D. Reward expectation selectively boosts the firing of
1268 accumbens D1+ neurons during motivated approach. *bioRxiv* 2023.09.02.556060 (2023)
1269 doi:10.1101/2023.09.02.556060.
- 1270 101. Martiros, N., Kapoor, V., Kim, S. E. & Murthy, V. N. Distinct representation of cue-
1271 outcome association by D1 and D2 neurons in the ventral striatum’s olfactory tubercle. *eLife*
1272 **11**, e75463 (2022).
- 1273 102. Nishioka, T. *et al.* Error-related signaling in nucleus accumbens D2 receptor-expressing
1274 neurons guides inhibition-based choice behavior in mice. *Nat. Commun.* **14**, 2284 (2023).
- 1275 103. Tsutsui-Kimura, I., Uchida, N. & Watabe-Uchida, M. Dynamical management of potential
1276 threats regulated by dopamine and direct- and indirect-pathway neurons in the tail of the
1277 striatum. *bioRxiv* 2022.02.05.479267 (2022) doi:10.1101/2022.02.05.479267.
- 1278 104. Gagnon, D. *et al.* Striatal Neurons Expressing D1 and D2 Receptors are Morphologically
1279 Distinct and Differently Affected by Dopamine Denervation in Mice. *Sci. Rep.* **7**, 41432
1280 (2017).
- 1281 105. Chen, R. *et al.* Decoding molecular and cellular heterogeneity of mouse nucleus accumbens.
1282 *Nat. Neurosci.* **24**, 1757–1771 (2021).
- 1283 106. Anderson, A. G., Kulkarni, A. & Konopka, G. A single-cell trajectory atlas of striatal

- 1284 development. *Sci. Rep.* **13**, 9031 (2023).
- 1285 107. Dana, H. *et al.* High-performance calcium sensors for imaging activity in neuronal
1286 populations and microcompartments. *Nat. Methods* **16**, 649–657 (2019).
- 1287 108. Gong, S. *et al.* A gene expression atlas of the central nervous system based on bacterial
1288 artificial chromosomes. *Nature* **425**, 917–925 (2003).
- 1289 109. Klapoetke, N. C. *et al.* Independent optical excitation of distinct neural populations. *Nat.*
1290 *Methods* **11**, 338–346 (2014).
- 1291 110. Govorunova, E. G., Sineshchekov, O. A., Janz, R., Liu, X. & Spudich, J. L.
1292 NEUROSCIENCE. Natural light-gated anion channels: A family of microbial rhodopsins
1293 for advanced optogenetics. *Science* **349**, 647–650 (2015).
- 1294 111. Li, N. *et al.* Spatiotemporal constraints on optogenetic inactivation in cortical circuits. *Elife*
1295 **8**, (2019).
- 1296 112. Lee, J. & Sabatini, B. L. Striatal indirect pathway mediates exploration via collicular
1297 competition. *Nature* **599**, 645–649 (2021).
- 1298 113. Collins, A. G. E. & Frank, M. J. Opponent actor learning (OpAL): modeling interactive
1299 effects of striatal dopamine on reinforcement learning and choice incentive. *Psychol. Rev.*
1300 **121**, 337–366 (2014).
- 1301 114. Gjorgjieva, J., Sompolinsky, H. & Meister, M. Benefits of pathway splitting in sensory
1302 coding. *J. Neurosci.* **34**, 12127–12144 (2014).
- 1303 115. Wang, Z., Wei, X.-X., Stocker, A. A. & Lee, D. D. Efficient Neural Codes under Metabolic
1304 Constraints. in *Advances in Neural Information Processing Systems* (eds. Lee, D.,
1305 Sugiyama, M., Luxburg, U., Guyon, I. & Garnett, R.) vol. 29 (Curran Associates, Inc.,
1306 2016).
- 1307 116. Ichinose, T. & Habib, S. ON and OFF Signaling Pathways in the Retina and the Visual
1308 System. *Front Ophthalmol (Lausanne)* **2**, (2022).
- 1309 117. Beier, K. T. *et al.* Circuit Architecture of VTA Dopamine Neurons Revealed by Systematic
1310 Input-Output Mapping. *Cell* **162**, 622–634 (2015).
- 1311 118. Poulin, J.-F., Gaertner, Z., Moreno-Ramos, O. A. & Awatramani, R. Classification of
1312 Midbrain Dopamine Neurons Using Single-Cell Gene Expression Profiling Approaches.
1313 *Trends Neurosci.* **43**, 155–169 (2020).
- 1314 119. Wenliang, L. K. *et al.* Distributional Bellman Operators over Mean Embeddings. *arXiv*
1315 *[stat.ML]* (2023).
- 1316 120. Niv, Y. Reinforcement learning in the brain. *J. Math. Psychol.* **53**, 139–154 (2009).
- 1317 121. Niv, Y., Edlund, J. A., Dayan, P. & O’Doherty, J. P. Neural prediction errors reveal a risk-
1318 sensitive reinforcement-learning process in the human brain. *J. Neurosci.* **32**, 551–562
1319 (2012).
- 1320 122. Necker, L. A. LXI. Observations on some remarkable optical phenomena seen in
1321 Switzerland; and on an optical phenomenon which occurs on viewing a figure of a crystal
1322 or geometrical solid. *The London, Edinburgh, and Dublin Philosophical Magazine and*
1323 *Journal of Science* **1**, 329–337 (1832).

- 1324 123. Gershman, S. J. & Uchida, N. Believing in dopamine. *Nat. Rev. Neurosci.* **20**, 703–714
1325 (2019).
- 1326 124. Wang, K., Zhou, K., Wu, R., Kallus, N. & Sun, W. The Benefits of Being Distributional:
1327 Small-Loss Bounds for Reinforcement Learning. *arXiv [cs.LG]* (2023).
- 1328 125. Luis, C. E., Bottero, A. G., Vinogradskaya, J., Berkenkamp, F. & Peters, J. Value-
1329 Distributional Model-Based Reinforcement Learning. *arXiv [cs.LG]* (2023).
- 1330 126. Chandak, Y. *et al.* Universal Off-Policy Evaluation. *arXiv [cs.LG]* (2021).
- 1331 127. Kim, D., Lee, K. & Oh, S. Trust Region-Based Safe Distributional Reinforcement Learning
1332 for Multiple Constraints. in *37th Conference on Neural Information Processing Systems*
1333 (2023).
- 1334 128. Kastner, T., Erdogdu, M. A. & Farahmand, A.-M. Distributional Model Equivalence for
1335 Risk-Sensitive Reinforcement Learning. *arXiv [cs.LG]* (2023).
- 1336 129. Cai, X.-Q. *et al.* Distributional Pareto-Optimal Multi-Objective Reinforcement Learning. in
1337 *37th Conference on Neural Information Processing Systems* (2023).
- 1338 130. Rigter, M., Lacerda, B. & Hawes, N. One Risk to Rule Them All: A Risk-Sensitive
1339 Perspective on Model-Based Offline Reinforcement Learning. *arXiv [cs.LG]* (2022).
- 1340 131. Bar-Gad, I., Morris, G. & Bergman, H. Information processing, dimensionality reduction
1341 and reinforcement learning in the basal ganglia. *Prog. Neurobiol.* **71**, 439–473 (2003).
- 1342 132. Barth-Maron, G. *et al.* Distributed Distributional Deterministic Policy Gradients. *arXiv*
1343 *[cs.LG]* (2018).
- 1344 133. Tessler, C., Efroni, Y. & Mannor, S. Action Robust Reinforcement Learning and
1345 Applications in Continuous Control. *arXiv [cs.LG]* (2019).
- 1346 134. Kuznetsov, A., Shvechikov, P., Grishin, A. & Vetrov, D. Controlling Overestimation Bias
1347 with Truncated Mixture of Continuous Distributional Quantile Critics. *arXiv [cs.LG]*
1348 (2020).
- 1349 135. Nam, D. W., Kim, Y. & Park, C. Y. GMAC: A Distributional Perspective on Actor-Critic
1350 Framework. in *Proceedings of the 38th International Conference on Machine Learning* (eds.
1351 Meila, M. & Zhang, T.) vol. 139 7927–7936 (PMLR, 18–24 Jul 2021).
- 1352 136. Duan, J. *et al.* Distributional Soft Actor-Critic: Off-Policy Reinforcement Learning for
1353 Addressing Value Estimation Errors. *IEEE Trans Neural Netw Learn Syst* **33**, 6584–6598
1354 (2022).
- 1355 137. O’Doherty, J. *et al.* Dissociable roles of ventral and dorsal striatum in instrumental
1356 conditioning. *Science* **304**, 452–454 (2004).
- 1357 138. Nasrallah, N. A. *et al.* Risk preference following adolescent alcohol use is associated with
1358 corrupted encoding of costs but not rewards by mesolimbic dopamine. *Proc. Natl. Acad. Sci.*
1359 *U. S. A.* **108**, 5466–5471 (2011).
- 1360 139. Tymula, A. *et al.* Dynamic prospect theory: Two core decision theories coexist in the
1361 gambling behavior of monkeys and humans. *Sci Adv* **9**, eade7972 (2023).
- 1362 140. Vrieze, E. *et al.* Reduced reward learning predicts outcome in major depressive disorder.
1363 *Biol. Psychiatry* **73**, 639–645 (2013).

- 1364 141. Brown, V. M. *et al.* Reinforcement Learning Disruptions in Individuals With Depression
1365 and Sensitivity to Symptom Change Following Cognitive Behavioral Therapy. *JAMA*
1366 *Psychiatry* **78**, 1113–1122 (2021).
- 1367 142. Gueguen, M. C. M., Schweitzer, E. M. & Konova, A. B. Computational theory-driven
1368 studies of reinforcement learning and decision-making in addiction: What have we learned?
1369 *Curr Opin Behav Sci* **38**, 40–48 (2021).
- 1370 143. Gong, S. *et al.* Targeting Cre recombinase to specific neuron populations with bacterial
1371 artificial chromosome constructs. *J. Neurosci.* **27**, 9817–9823 (2007).
- 1372 144. Gerfen, C. R., Paletzki, R. & Heintz, N. GENSAT BAC cre-recombinase driver lines to
1373 study the functional organization of cerebral cortical and basal ganglia circuits. *Neuron* **80**,
1374 1368–1383 (2013).
- 1375 145. Thiele, S. L., Warre, R. & Nash, J. E. Development of a unilaterally-lesioned 6-OHDA
1376 mouse model of Parkinson’s disease. *J. Vis. Exp.* (2012) doi:10.3791/3234.
- 1377 146. Uchida, N. & Mainen, Z. F. Speed and accuracy of olfactory discrimination in the rat. *Nat.*
1378 *Neurosci.* **6**, 1224–1229 (2003).
- 1379 147. Jun, J. J. *et al.* Fully integrated silicon probes for high-density recording of neural activity.
1380 *Nature* **551**, 232–236 (2017).
- 1381 148. Steinmetz, N. A. *et al.* Neuropixels 2.0: A miniaturized high-density probe for stable, long-
1382 term brain recordings. *Science* **372**, (2021).
- 1383 149. Pachitariu, M., Sridhar, S. & Stringer, C. Solving the spike sorting problem with Kilosort.
1384 *bioRxiv* 2023.01.07.523036 (2023) doi:10.1101/2023.01.07.523036.
- 1385 150. Zhou, Z. C. *et al.* Deep-brain optical recording of neural dynamics during behavior. *Neuron*
1386 **111**, 3716–3738 (2023).
- 1387 151. Pachitariu, M. *et al.* Suite2p: Beyond 10,000 neurons with standard two-photon microscopy.
1388 (2016) doi:10.1101/061507.
- 1389 152. Friedrich, J., Zhou, P. & Paninski, L. Fast online deconvolution of calcium imaging data.
1390 *PLoS Comput. Biol.* **13**, e1005423 (2017).
- 1391 153. Lopes, G. *et al.* Bonsai: an event-based framework for processing and controlling data
1392 streams. *Front. Neuroinform.* **9**, 7 (2015).
- 1393 154. Pisanello, M. *et al.* Tailoring light delivery for optogenetics by modal demultiplexing in
1394 tapered optical fibers. *Sci. Rep.* **8**, 4467 (2018).
- 1395 155. Lee, J., Wang, W. & Sabatini, B. L. Anatomically segregated basal ganglia pathways allow
1396 parallel behavioral modulation. *Nat. Neurosci.* **23**, 1388–1398 (2020).
- 1397 156. Sanders, J. I. & Kepecs, A. A low-cost programmable pulse generator for physiology and
1398 behavior. *Front. Neuroeng.* **7**, 43 (2014).
- 1399 157. Shamash, P., Carandini, M., Harris, K. & Steinmetz, N. A tool for analyzing electrode tracks
1400 from slice histology. *bioRxiv* 447995 (2018) doi:10.1101/447995.
- 1401 158. Wang, Q. *et al.* The Allen Mouse Brain Common Coordinate Framework: A 3D Reference
1402 Atlas. *Cell* **181**, 936-953.e20 (2020).
- 1403 159. Claudi, F. *et al.* Visualizing anatomically registered data with brainrender. *eLife* **10**, (2021).

- 1404 160. Chon, U., Vanselow, D. J., Cheng, K. C. & Kim, Y. Enhanced and unified anatomical
1405 labeling for a common mouse brain atlas. *Nat. Commun.* **10**, 5067 (2019).
- 1406 161. Claudi, F. *et al.* BrainGlobe Atlas API: a common interface for neuroanatomical atlases. *J.*
1407 *Open Source Softw.* **5**, 2668 (2020).
- 1408 162. Franklin, K. B. J. & Paxinos, G. *Paxinos and Franklin's The mouse brain in stereotaxic*
1409 *coordinates*. (Academic Press, an imprint of Elsevier, 2013).
- 1410 163. Hintiryan, H. *et al.* The mouse cortico-striatal projectome. *Nat. Neurosci.* **19**, 1100–1114
1411 (2016).
- 1412 164. Peters, A. J., Fabre, J. M. J., Steinmetz, N. A., Harris, K. D. & Carandini, M. Striatal activity
1413 topographically reflects cortical activity. *Nature* **591**, 420–425 (2021).
- 1414 165. Harris, C. R. *et al.* Array programming with NumPy. *Nature* **585**, 357–362 (2020).
- 1415 166. Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in Python.
1416 *Nat. Methods* **17**, 261–272 (2020).
- 1417 167. McKinney, W. Data Structures for Statistical Computing in Python. in *Proceedings of the*
1418 *9th Python in Science Conference* (eds. van der Walt, S. & Millman, J.) (2010).
1419 doi:10.25080/majora-92bf1922-00a.
- 1420 168. Buitinck, L. *et al.* API design for machine learning software: experiences from the scikit-
1421 learn project. in *ECML PKDD Workshop: Languages for Data Mining and Machine*
1422 *Learning* 108–122 (2013).
- 1423 169. Seabold, S. & Perktold, J. Statsmodels: Econometric and statistical modeling with python. in
1424 *Proceedings of the 9th Python in Science Conference* (SciPy, 2010). doi:10.25080/majora-
1425 92bf1922-011.
- 1426 170. Hunter, J. D. Matplotlib: A 2D Graphics Environment. *Comput. Sci. Eng.* **9**, 90–95 (May-
1427 June 2007).
- 1428 171. Waskom, M. seaborn: statistical data visualization. *J. Open Source Softw.* **6**, 3021 (2021).
- 1429 172. Dietterich, T. G. Approximate Statistical Tests for Comparing Supervised Classification
1430 Learning Algorithms. *Neural Comput.* **10**, 1895–1923 (1998).
- 1431 173. Rescorla, R. A. & Wagner, A. R. A Theory of Pavlovian Conditioning: Variations in the
1432 Effectiveness of Reinforcement and Nonreinforcement. in *Classical conditioning II: current*
1433 *research and theory* (ed. A H Black & W) 64–99 (Appleton-Century-Crofts, 1972).
- 1434 174. Huber, P. J. Robust Estimation of a Location Parameter. *Ann. Math. Stat.* **35**, 73–101
1435 (1964).
- 1436 175. Romero Pinto, S. & Uchida, N. Tonic dopamine and biases in value learning linked through
1437 a biologically inspired reinforcement learning model. *bioRxiv* 2023.11.10.566580 (2023)
1438 doi:10.1101/2023.11.10.566580.
- 1439 176. Gagne, C. & Dayan, P. Peril, prudence and planning as risk, avoidance and worry. *J. Math.*
1440 *Psychol.* **106**, 102617 (2022).
- 1441 177. Rockafellar, R. T. & Uryasev, S. Optimization of conditional value-at-risk. *Journal of Risk*
1442 **2**, 21–41 (2000).
- 1443 178. Churchland, M. M. *et al.* Stimulus onset quenches neural variability: a widespread cortical

1444 phenomenon. *Nat. Neurosci.* **13**, 369–378 (2010).

1445 **Acknowledgments**

1446 We thank members of the Uchida Lab for valuable discussions and comments on the manuscript.
1447 Ed Soucy and Brett Graham of the Harvard Center for Brain Science Neurotechnology Core
1448 Facility provided critical assistance with instrumentation. We'd also like to thank Dr. Allison
1449 Girasole and Prof. Bernardo Sabatini for sharing the GtACR1 mouse line; Dr. Xintong Cai, Prof.
1450 Bernardo Sabatini, Prof. Chris Harvey, and Prof. Sam Gershman for helpful conversations; and
1451 Dr. Matteo Carandini, Dr. Kenneth Harris, Dr. Andrew Peters and other members of the Cortex
1452 lab for their advice on Neuropixels recording. This work was supported by grants from NIH
1453 (R01NS116753, to N.U. and J.D.; F31NS124095, to A.S.L.), the Human Frontier Science
1454 Program (LT000801/2018, to S.M.), the Harvard Brain Science Initiative, and the Brain and
1455 Behavior Research Foundation (NARSAD Young Investigator no. 30035 to S.M.). We thank the
1456 Harvard Center for Biological Imaging (RRID:SCR_018673) for infrastructure and support for
1457 *ex vivo* imaging, which was funded in part by the Simmons Award (to A.S.L.). The computations
1458 in this paper were run in part on the FASRC Cannon cluster supported by the FAS Division of
1459 Science Research Computing Group at Harvard University.

1460 **Author Contributions**

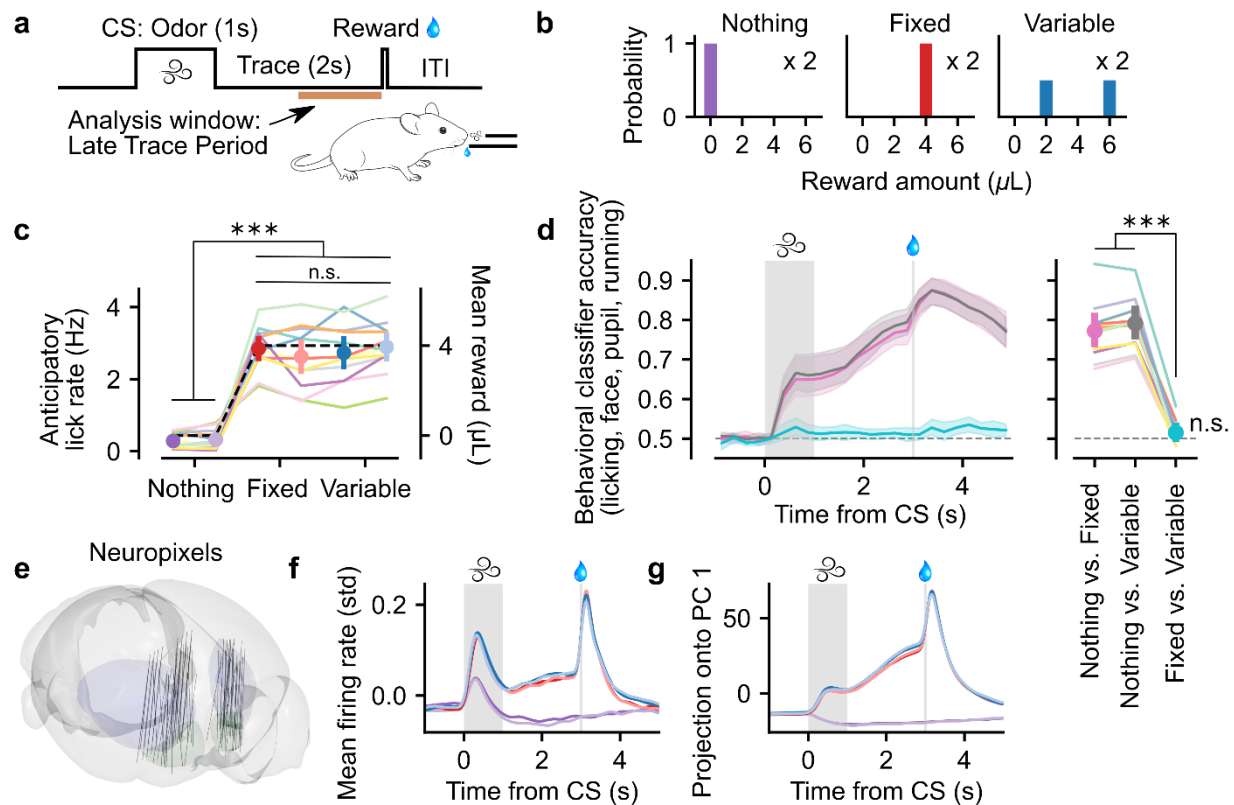
1461 A.S.L. and N.U. designed the experiments. A.S.L. and M.M. performed the experiments, with
1462 initial help from S.M. A.S.L. and M.M. preprocessed the data. A.S.L. analyzed the data and
1463 implemented the computational models with input from J.D. and N.U. Q.Z. implemented ANN-
1464 based distributional decoding under the supervision of J.D. A.S.L. wrote the first draft of the
1465 manuscript and created the figures. N.U., J.D., S.M., and A.S.L. edited the manuscript.

1466 **Competing Interests**

1467 The authors declare no conflicts of interest.

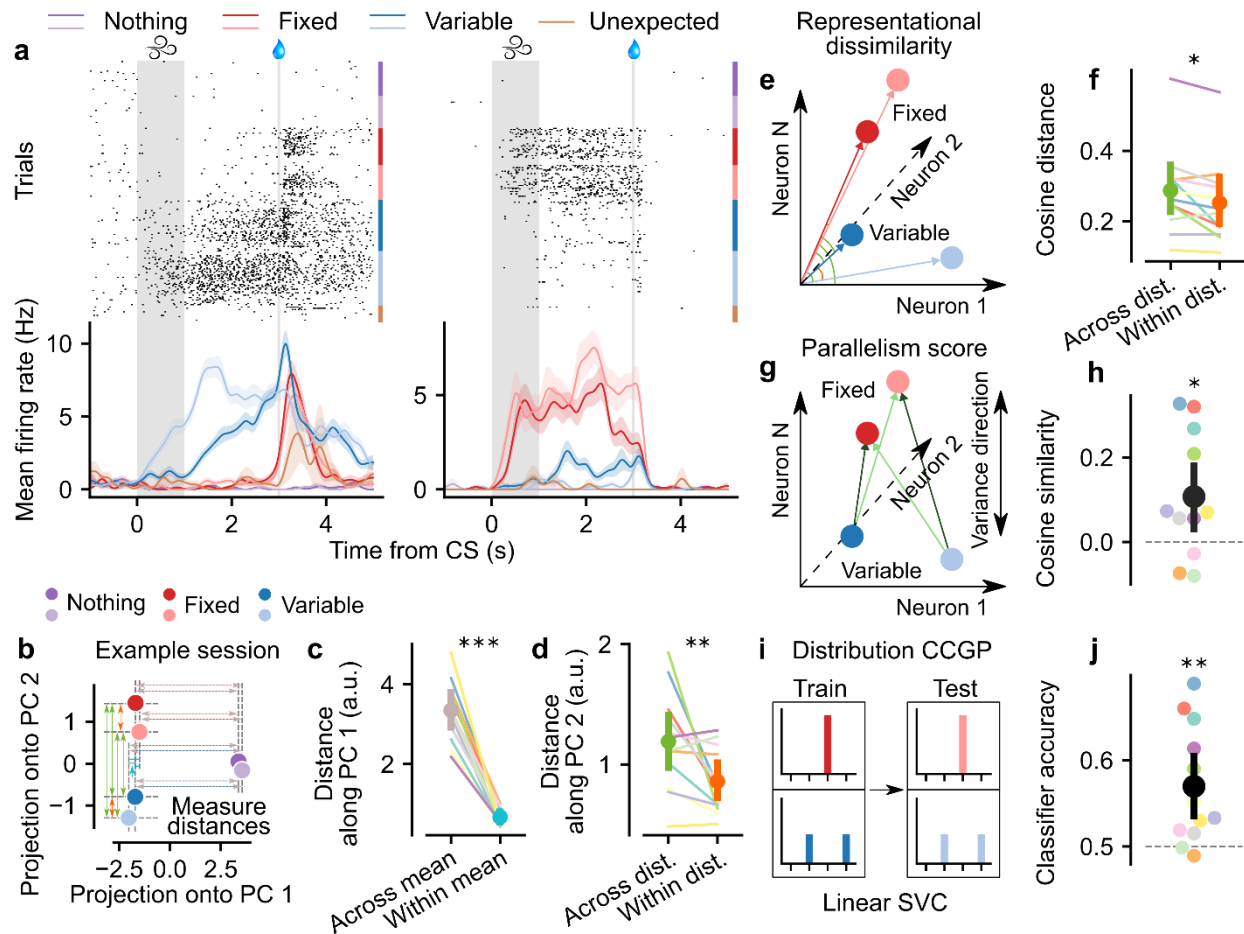
1468 **Materials and Correspondence**

1469 Please direct any requests for materials to Naoshige Uchida.



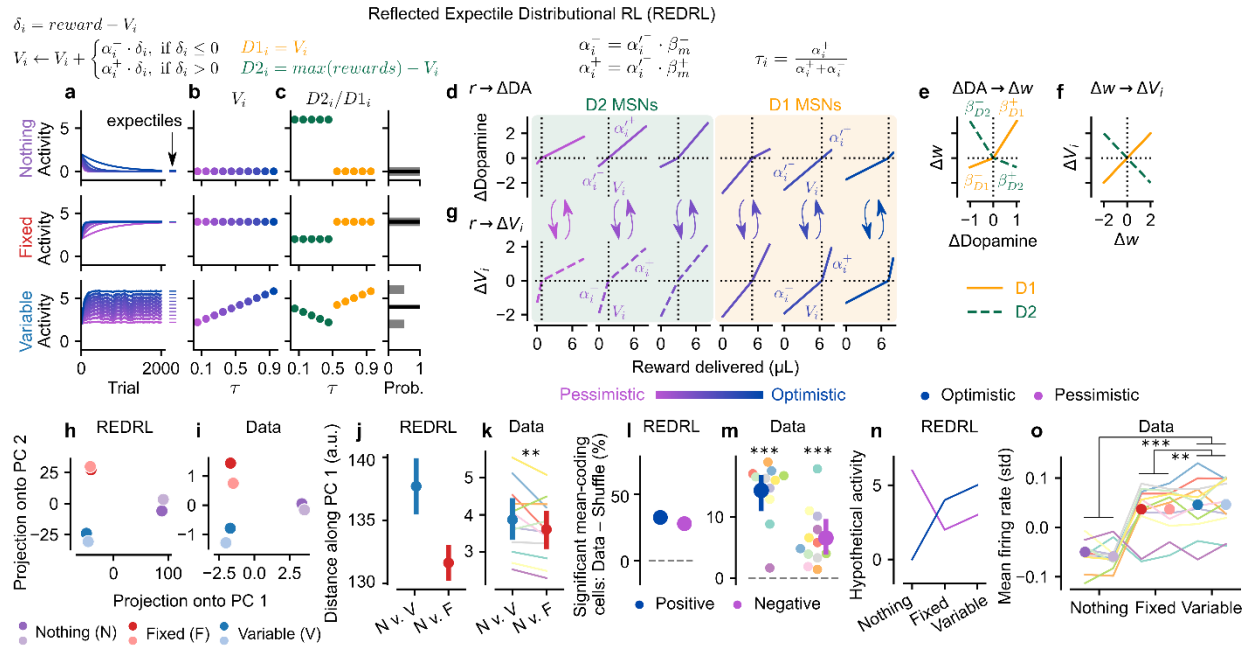
1470 **Fig. 1 | A classical conditioning task and recording setup to investigate distributional**
 1471 **reinforcement learning.** **a**, Water-restricted, head-fixed mice were trained to associate odors
 1472 with stochastic rewards following a brief (2 s) trace period. When not otherwise specified,
 1473 behavioral and neural activity were analyzed in the final second of the trace period (“Late Trace”
 1474 period) in order to assess reward anticipation. Odor-reward distribution mappings were
 1475 randomized across mice. CS, conditioned stimulus. ITI, inter-trial interval. **b**, Probability
 1476 distributions over reward amounts that were paired with odors. Each distribution was associated
 1477 with two distinct odors, for a total of six odors, in order to distinguish stimulus information from
 1478 distributional content. Furthermore, two distributions (Fixed and Variable) had the same mean of
 1479 4 μL , but different variance. **c**, Anticipatory lick rates for each trial type, computed during the
 1480 Late Trace period (Nothing 1 or Nothing 2: $p < 0.001$ versus Fixed 1, Fixed 2, Variable 1, and
 1481 Variable 2; Fixed 1: $p = 0.502, 0.925, 0.419$ versus Fixed 2, Variable 1, and Variable 2,
 1482 respectively). **d**, Cross-validated classification accuracy of a linear kernel Support Vector
 1483 Machine trained on licking, pupil area, whisking, running, and singular value decomposition of
 1484 behavioral videos (Extended Data Fig. 1). The data associated with the two odors corresponding
 1485 to the same distribution were pooled and then split into training and validation sets. *Left*,
 1486 behavioral classifier accuracy across time. Predictors were aggregated within 250 ms, non-
 1487 overlapping bins. Shaded regions denote 95% confidence intervals across mice. Pink, Nothing
 1488 vs. Fixed; Grey, Nothing vs. Variable; Cyan, Fixed vs. Variable. *Right*, quantification of
 1489 behavioral classifier accuracy when trained separately on the entire Late Trace period (Fixed vs.

1490 Variable: $p < 0.001$ versus Nothing vs. Fixed and Nothing vs. Variable, $p = 0.053$ compared to
1491 chance level of 50%). **e**, Reconstructed Neuropixels probe trajectories, aligned to the Allen
1492 Mouse Brain Common Coordinate Framework. **f**, Individual neurons' firing rates were z-scored
1493 across time, aligned to stimulus onset, averaged for each trial type, and then averaged across
1494 neurons. Color code as in **c**. Average firing rates correlate with mean reward. **g**, Trial type
1495 averages for each neuron were concatenated, and the first principal component was extracted and
1496 plotted across neurons. Color code as in **c**. For Figs. 1–3 and 6, asterisks represent the result of
1497 Linear Mixed Effects model across sessions with a random intercept for each mouse, and, if
1498 applicable, a random slope for each mouse as a function of grouping (e.g. Across- vs. Within-
1499 distribution): ***, $p < 0.001$; **, $p < 0.01$; *, $p < 0.05$; n.s., not significant at $\alpha = 0.05$. Asterisks
1500 over lines connecting different groupings indicate significant differences between groups, while
1501 asterisks without corresponding lines indicate that the group is significantly different from
1502 chance, indicated by the dashed grey line. The shaded region from 0 to 1 s represents the interval
1503 of odor delivery, and the vertical line at 3 s indicates reward timing. For Figures 1–4, pastel
1504 colors in the background show averages across sessions within mice, while dots with error bars
1505 in the foreground denote means and 95% confidence intervals across mice. Differences were
1506 taken within-session.



1507 **Fig. 2 | Distributional coding across the striatum.** **a**, Example peri-stimulus time histograms
 1508 (PSTHs) of two simultaneously-recorded neurons in the ventromedial striatum. *Top*, spike
 1509 rasters, aligned to odor onset and sorted by trial type. *Bottom*, mean \pm s.e.m. firing to each trial
 1510 type, after smoothing the entire session's spike train with a Gaussian kernel (s.d. = 100 ms).
 1511 While both neurons tend to increase on average to rewarded odors, the neuron on the left prefers
 1512 Variable odors, while the one on the right prefers Fixed odors, and tend to do so consistently for
 1513 different odors associated with the same distribution. **b**, Firing rate during the Late Trace period,
 1514 averaged across trials of each type, was projected into two dimensions using principal
 1515 component analysis (PCA) independently for each session. We then measured the distances
 1516 between trial types along each PC, as shown by the arrows. Color code as in **a**. **c**, Euclidean
 1517 distance along PC 1 was significantly greater for across-mean pairs (Nothing vs. Rewarded) than
 1518 within-mean pairs (Fixed vs. Variable; $p < 0.001$). **d**, Euclidean distance along PC 2 was
 1519 significantly greater for across-distribution pairs (Fixed vs. Variable) than within-distribution
 1520 pairs (Fixed 1 vs. Fixed 2 or Variable 1 vs. Variable 2; $p = 0.006$). **e**, Schematic illustrating
 1521 representational dissimilarity analysis (RDA). The population vector corresponding to each trial
 1522 type was computed independently for each session. We then computed the cosine distances
 1523 between across-distribution and within-distribution pairs, shown by the green and orange arcs. **f**,
 1524 Quantification of cosine distances (Across- vs. Within-distribution: $p = 0.029$). **g**, Schematic

1525 illustrating parallelism score. We computed the difference vector between each Fixed and
1526 Variable trial type for each session independently. Parallelism score is defined as the cosine
1527 similarity between each non-overlapping pair of vectors, averaged over the two possible
1528 combinations (dark green and light green). **h**, Quantification of parallelism score ($p = 0.015$
1529 compared to chance level of 0). **i**, Schematic illustrating computation of cross-condition
1530 generalization performance (CCGP). Linear support vector classifiers (SVCs) were trained to
1531 discriminate one Fixed and one Variable odor and then tested on the held-out Fixed vs. Variable
1532 pair. This was then repeated and averaged over all four possible combinations of training and test
1533 sets. **j**, Quantification of CCGP ($p = 0.001$ compared to chance level of 50%).



1534

1535

Fig. 3 | Reflected Expectile Distributional Reinforcement Learning (REDRL). a-c,

1536

Algorithmic REDRL model. **a**, Over the course of training, value predictors (V_i , here initialized

1537

to 2) converge to the expectiles of the associated reward distribution. **b**, Post-learning activity of

1538

the simulated value predictors, V_i , as a function of their optimism level. The relative pessimism

1539

or optimism of each predictor is parameterized by τ , which can range from 0 to 1 (x -axis). **c**, *Left*,

1540

pessimistic ($\tau < 0.5$) value predictors are identified with D2 MSNs (green), and their coding is

1541

flipped such that decreases in D2 activity correspond to increases in V_i , and vice versa.

1542

Optimistic ($\tau > 0.5$) predictors are directly proportional to D1 MSN activity (orange). *Right*,

1543

collectively, this striatal code characterizes the complete reward distribution via its expectiles. **d-g**,

1544

Implementation of REDRL within the mesolimbic circuit. **d**, Heterogeneity across dopamine

1545

neurons can be characterized using piecewise linear functions. Pessimistic neurons have high

1546

slopes in the negative domain (α_i^-) and low slopes in the positive domain (α_i^+), while the

1547

opposite is true for optimistic neurons. Over the course of learning, the zero-crossing point V_i

1548

associated with each neuron will shift to equal the τ_i -th expectile (vertical dotted line)⁴. **e**, D1 and

1549

D2 MSNs have asymmetric plasticity rules, potentiating more to increases and decreases in

1550

dopamine, respectively, relative to baseline (vertical dotted line)³⁹. **f**, As a consequence, D1

1551

activity is expected to correlate positively, and D2 activity negatively, with the corresponding

1552

value prediction^{5,6}. To recover V_i , we must subtract out the D2 activity, which could be

1553

accomplished for instance via its inhibitory projection to the ventral pallidum. **g**, The change in

1554

each value predictor is $\Delta V_i = \alpha_i^{-/+} \cdot \beta_m^{-/+} \cdot \delta = \alpha_i^{-/+} \cdot \delta$, as demanded by the gradient

1555

descent-based update rule. The net result is that D1 MSNs are biased optimistically, and D2

1556

pessimistically, relative to their dopamine input asymmetries, because their learning constants

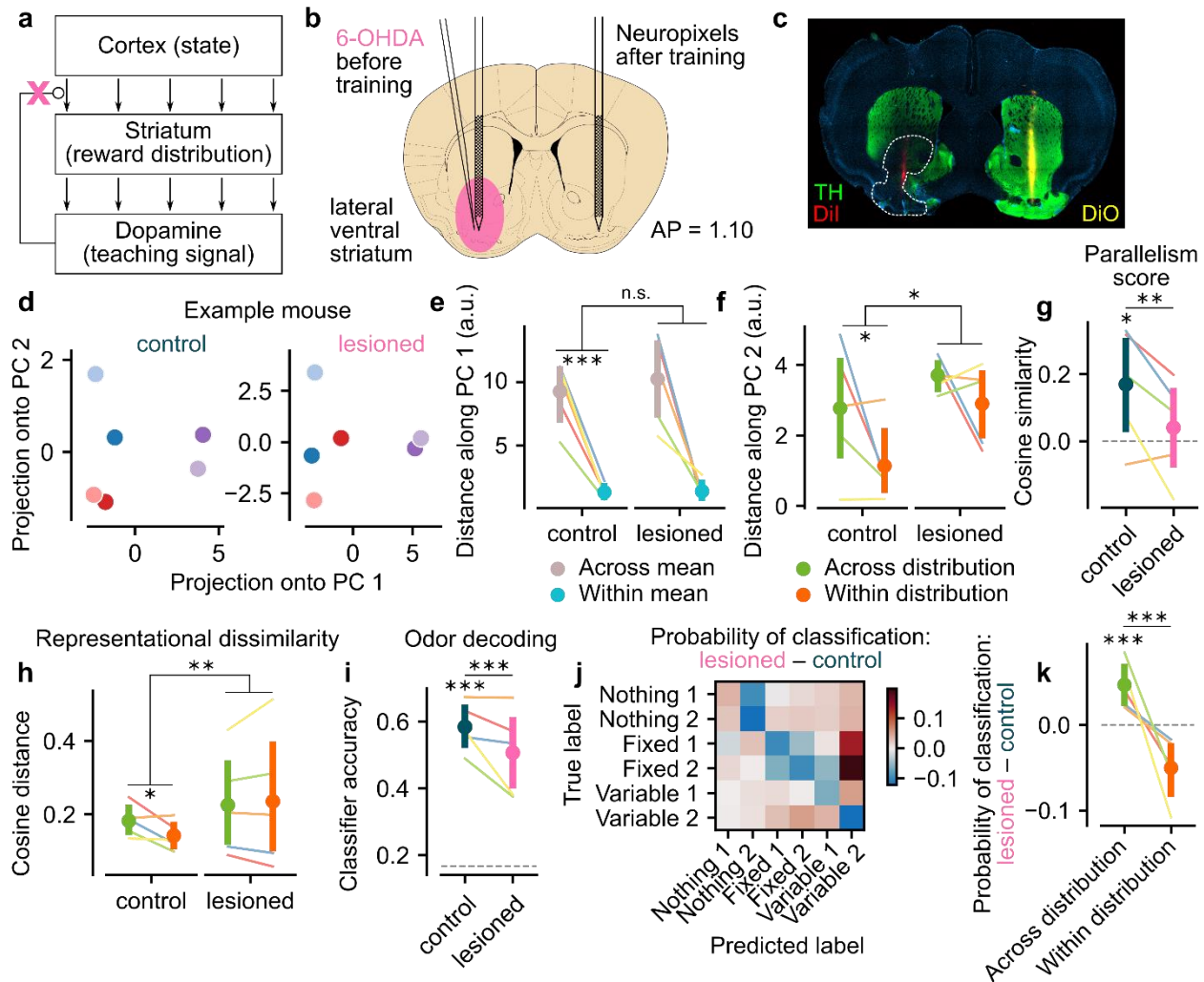
1557

$\beta_m^{-/+}$ favor positive and negative prediction errors, respectively. With the plasticity rule shown,

1558

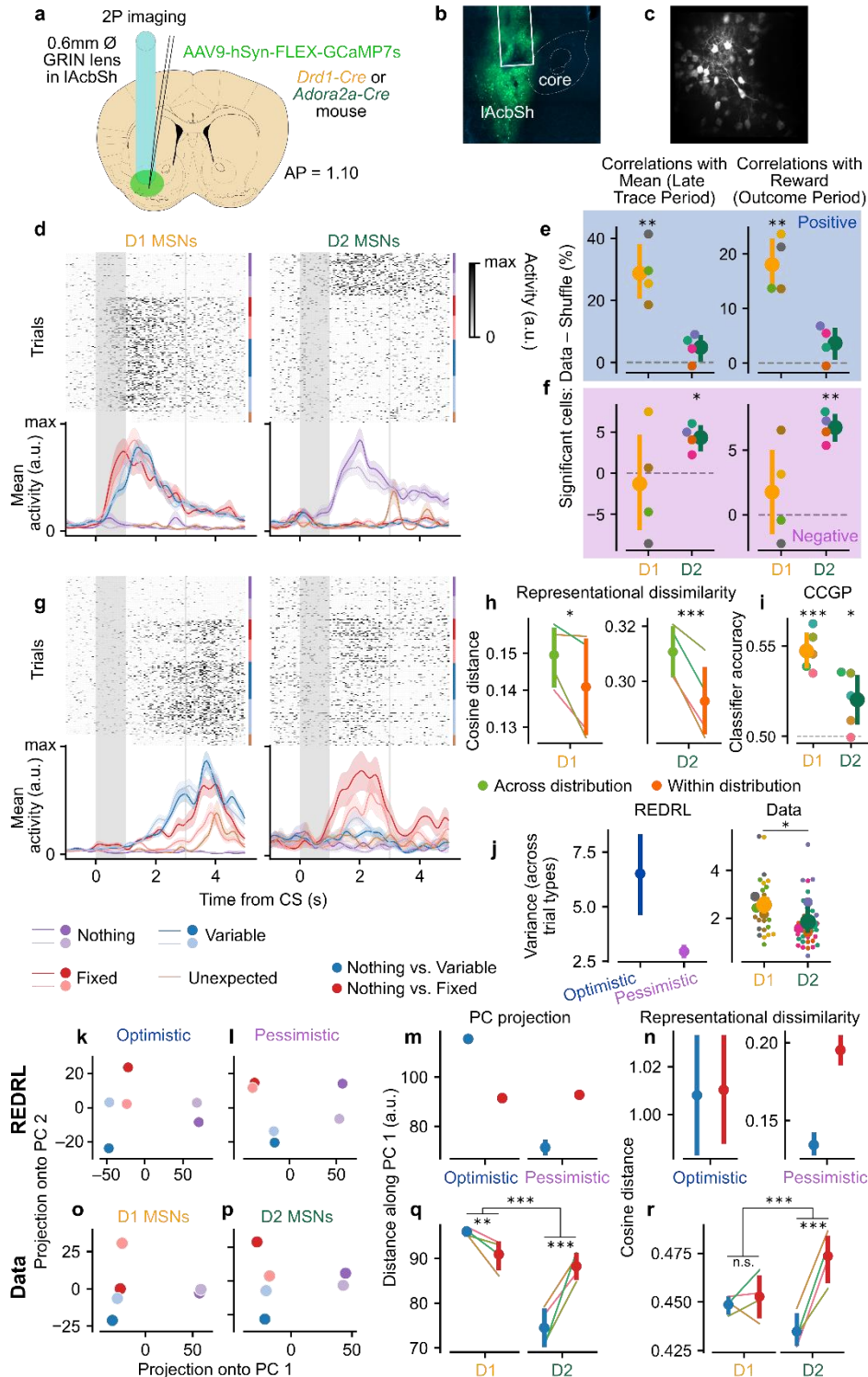
all D1 MSNs have $\tau > 0.5$, and all D2 MSNs have $\tau < 0.5$, justifying the division in **c**, though the

1559 precise distribution will depend on the specific plasticity rule and distribution of dopamine
1560 asymmetries. **h**, Two-dimensional PCA projection of converged value predictors, plus noise, for
1561 the REDRL model. Variable odors are further separated than Fixed from Nothing along PC 1
1562 because after mean-centering, the patterns of Nothing and Variable activity are almost perfectly
1563 anticorrelated with one another, and the PC 1 loadings closely resemble Nothing activity itself. **i**,
1564 PCA projection of example session (same as Fig. 2b) shows a striking resemblance to the
1565 REDRL prediction in separating primarily Rewarded and Unrewarded odors along PC 1 and
1566 Fixed and Variable odors along PC 2. **j**, In addition, REDRL predicts that the distance between
1567 Nothing (N) and Variable (V) odors along PC 1 should be slightly greater than that between
1568 Nothing and Fixed (F). **k**, Striatal data are consistent with this prediction, with the distance along
1569 PC 1 significantly greater for Nothing vs. Variable than Nothing vs. Fixed odor pairs ($p = 0.007$).
1570 **l**, REDRL predicts that there should be substantial fractions of neurons that correlate either
1571 positively or negatively with mean value, corresponding to D1 and D2 MSNs. **m**, Significant
1572 populations of striatal neurons encode mean reward positively and negatively. Mean reward
1573 predicted on each trial was correlated with Late Trace activity. Then, for each neuron
1574 independently, we shuffled the odor-distribution mappings and re-computed the correlations.
1575 Each point denotes the per-mouse difference in fraction of significant cells (that is, cells with
1576 uncorrected $p < 0.05$) for the unshuffled and shuffled data, separately for cells that correlated
1577 positively or negatively with mean reward (Positive and Negative: $p < 0.001$, paired samples t -
1578 test comparing ordered and shuffled fractions across mice). **n**, REDRL predicts that Variable
1579 odors elicit higher population mean firing than Fixed odors, regardless of the optimism or
1580 pessimism of the underlying value predictor. **o**, Mean z-scored firing rates for each neuron, in
1581 addition to being higher for Rewarded than Unrewarded odors ($p < 0.001$), were also higher for
1582 Variable than for Fixed odors ($p = 0.006$), as assessed by an LME with neuron-level
1583 observations, averaged over trials, and mouse-level random effects.



1584 **Fig. 4 | Dopamine is necessary for learning distributional representations.** **a**, Schematic
 1585 illustration of the basal ganglia, showing how dopamine is hypothesized to act as a teaching
 1586 signal to update corticostriatal synaptic weights. Therefore, dopamine lesions (pink “x”) are
 1587 predicted to disrupt representations of the reward distribution in the striatum. **b**, Schematic
 1588 illustration of dopamine lesion experiment. 6-OHDA was injected unilaterally into the lateral
 1589 ventral striatum of naive mice to ablate dopamine neurons. After recovery and training, we
 1590 recorded striatal activity in both the lesioned and control hemispheres. **c**, Histology from an
 1591 example 6-OHDA animal showing Neuropixels probe tracks (red and yellow), dopamine axons
 1592 (green), and lesion (white dashed line surrounding region of reduced TH staining). **d**, PCA
 1593 projection of Late Trace activity from the control (*left*) and lesioned (*right*) hemispheres for an
 1594 example mouse. **e**, Distance along PC 1, while significantly higher for across-mean than within-
 1595 mean pairs ($p < 0.001$), does not differ between hemispheres ($p = 0.676$). For all panels of this
 1596 figure, colored lines denote individual mice, averaged across pseudo-populations, and LMEs use
 1597 these pseudo-populations as the individual observations with mouse-level random effects. **f**, By
 1598 contrast, the difference in distance along PC 2 between across- and within-distribution pairs is
 1599 significantly positive ($p = 0.033$) and greater for the control relative to the lesioned hemisphere

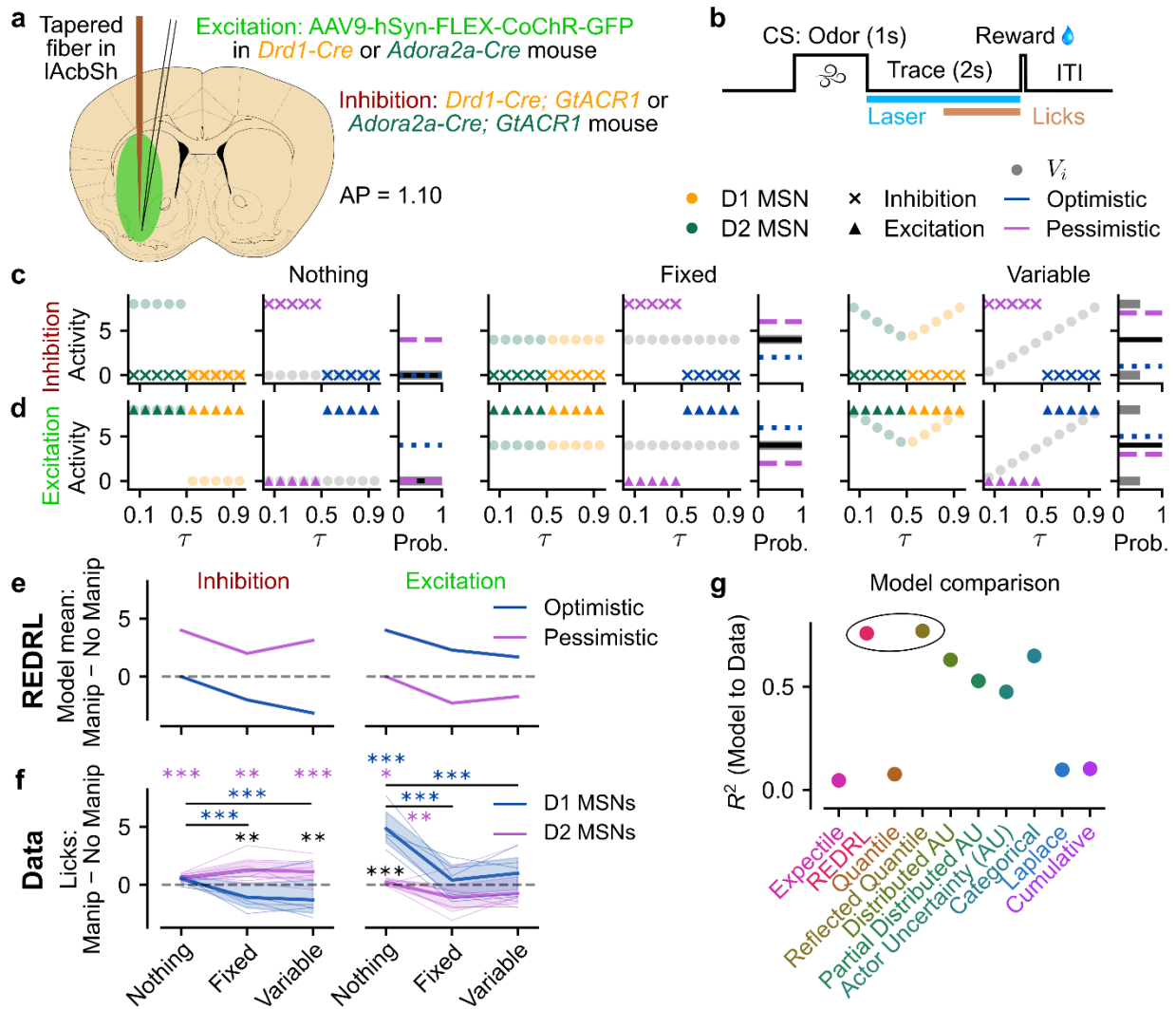
1600 ($p = 0.026$). **g**, Parallelism score is significantly positive ($p = 0.029$) and greater in the control
1601 relative to the lesioned hemisphere ($p = 0.009$). **h**, Similarly, the difference in representational
1602 dissimilarity between across- and within-distribution pairs is significantly positive ($p = 0.036$)
1603 and greater in the control relative to the lesioned hemisphere ($p = 0.005$). **i**, Six-way odor
1604 classification accuracy during the Odor period is above chance ($p < 0.001$) and higher for the
1605 control relative to the lesioned hemisphere ($p < 0.001$). **j**, Difference in odor classifier confusion
1606 matrices between the control and lesioned hemispheres. The probability of correct classification
1607 (main diagonal) decreases for nearly all trial types upon lesioning. **k**, The decrement in odor
1608 coding due to the lesion is mainly due to an increase in across-distribution, within-mean
1609 classification errors (the tendency in the lesioned hemisphere to predict Variable even when the
1610 true label was Fixed; $p < 0.001$) and a concomitant decrease in within-distribution classification
1611 ($p < 0.001$ for Across- vs. Within-distribution difference).



1612

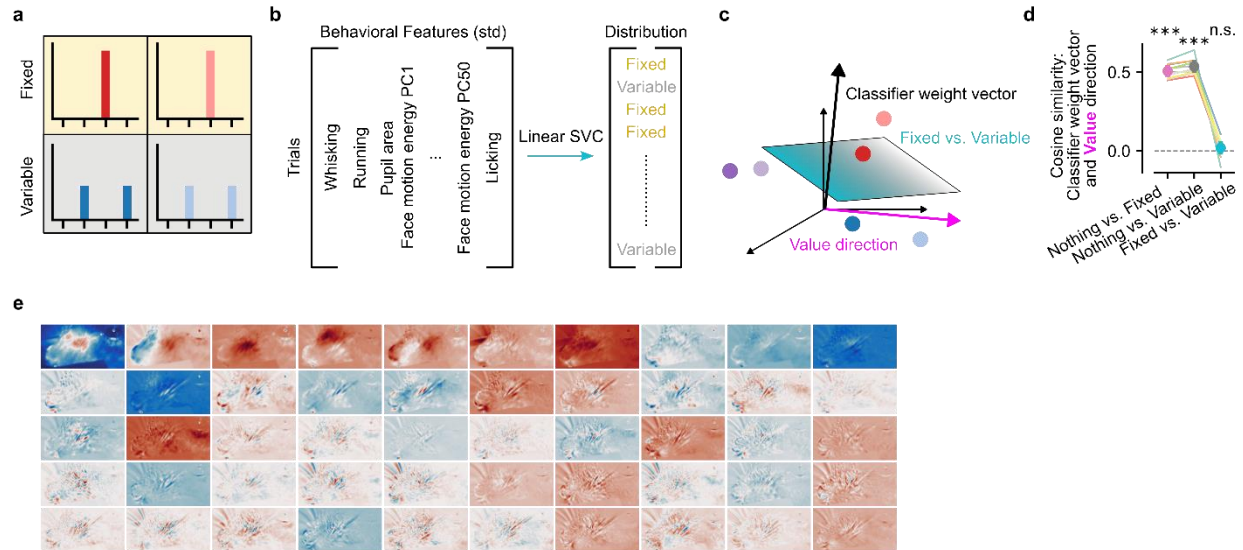
1613 **Fig. 5 | Opponent contributions of D1 and D2 MSNs to distributional coding.** **a**, Schematic
1614 illustration of two-photon calcium imaging experiment. We first injected a virus encoding the
1615 calcium indicator GCaMP7s and then implanted a GRIN lens in the lAcbSh in either *Drd1-Cre*
1616 or *Adora2a-Cre* mice, which drive Cre-dependent expression specifically in D1 and D2 MSNs,

1617 respectively. **b**, Example slice showing expression of GCaMP in the lAcbSh in a *Drd1-Cre*
1618 animal. **c**, Example FOV imaged through a GRIN lens in a *Drd1-Cre* animal. **d**, Deconvolved
1619 Ca^{2+} activity from an example D1 (*left*) and D2 (*right*) MSN. As in Fig. 2a, the top panel is a
1620 raster plot, normalized by maximum deconvolved activity, and the bottom panel shows average
1621 deconvolved activity \pm s.e.m. across trials of each type. The D1 MSN responds most to
1622 Rewarded odors, while the D2 MSN responds most to Nothing odors. **e**, Quantification of
1623 average percentage of cells that correlate significantly positively with mean (*left*) or reward
1624 (*right*) during the Late Trace and Outcome periods, respectively, relative to the expectation from
1625 odor coding alone (shuffling odor-distribution mappings, horizontal dashed line). There are
1626 significantly more cells than expected by chance for D1 (paired samples *t*-test comparing ordered
1627 and shuffled fractions across mice, $p = 0.009$, 0.006 for mean and reward, respectively), but not
1628 D2 ($p = 0.113$, 0.107) MSNs. Thick lines show the mean \pm 95% confidence interval across mice.
1629 **f**, Same as **e**, but for significant negative correlations. In this case, D2 ($p = 0.013$, 0.001) but not
1630 D1 ($p = 0.736$, 0.433) cells are significantly more common than expected by chance. **g**, Same as
1631 **d**, but showing an example D1 (*left*) and D2 (*right*) MSN that reliably discriminate Fixed and
1632 Variable odors. **h**, Cosine distance is significantly greater for across than within-distribution
1633 pairs for both D1 ($p = 0.022$) and D2 ($p < 0.001$) MSNs. For panels **h**, **i**, **q**, and **r** of this figure,
1634 individual replicates are pseudo-populations, split across trials and pooled across mice, hence
1635 there are no mouse-level random effects. Thick lines show the mean \pm 95% confidence interval
1636 across pseudo-populations. **i**, CCGP is significantly above chance for both D1 (one-sample *t*-test
1637 relative to 0.5, $p < 0.001$) and D2 ($p = 0.048$) MSNs, demonstrating abstract encoding of
1638 variance in both populations. **j**, Variance across trial types, computed for the simulated REDRL
1639 predictors (*left*) and normalized neural data (*right*). Small dots are averages within sessions,
1640 medium dots are averages within mice, and large dots with error bars show averages \pm 95%
1641 confidence intervals across mice ($p = 0.017$ for effect of genotype). **k-l**, Simulated REDRL value
1642 predictors were projected into two-dimensional PC space separately for optimistic (D1, **k**) or
1643 pessimistic (D2, **l**) value predictors. **m**, Quantification of Euclidean distance along PC 1 for the
1644 REDRL model. While optimistic predictors show the same trend as the complete code (Fig. 3j),
1645 pessimistic predictors swap the ordering between Fixed and Variable odors. Error bars denote
1646 95% confidence intervals across odor pairs. **n**, Same as **m**, but using cosine distance in the full-
1647 dimensional space to quantify representational dissimilarity, again independently for optimistic
1648 and pessimistic predictors. **o-r**, Same as **k-n**, but showing data collected from D1 and D2 MSNs
1649 rather than simulated optimistic and pessimistic predictors, respectively. For both the distance
1650 along PC 1 (Nothing vs. Variable compared to Nothing vs. Fixed: $p = 0.001$ for D1, $p < 0.001$
1651 for D2, $p < 0.001$ for the relative differences between D1 and D2) and the representational
1652 dissimilarity ($p = 0.489$ for D1, $p < 0.001$ for D2, $p < 0.001$ for the relative differences), striatal
1653 data closely match the theoretical predictions.

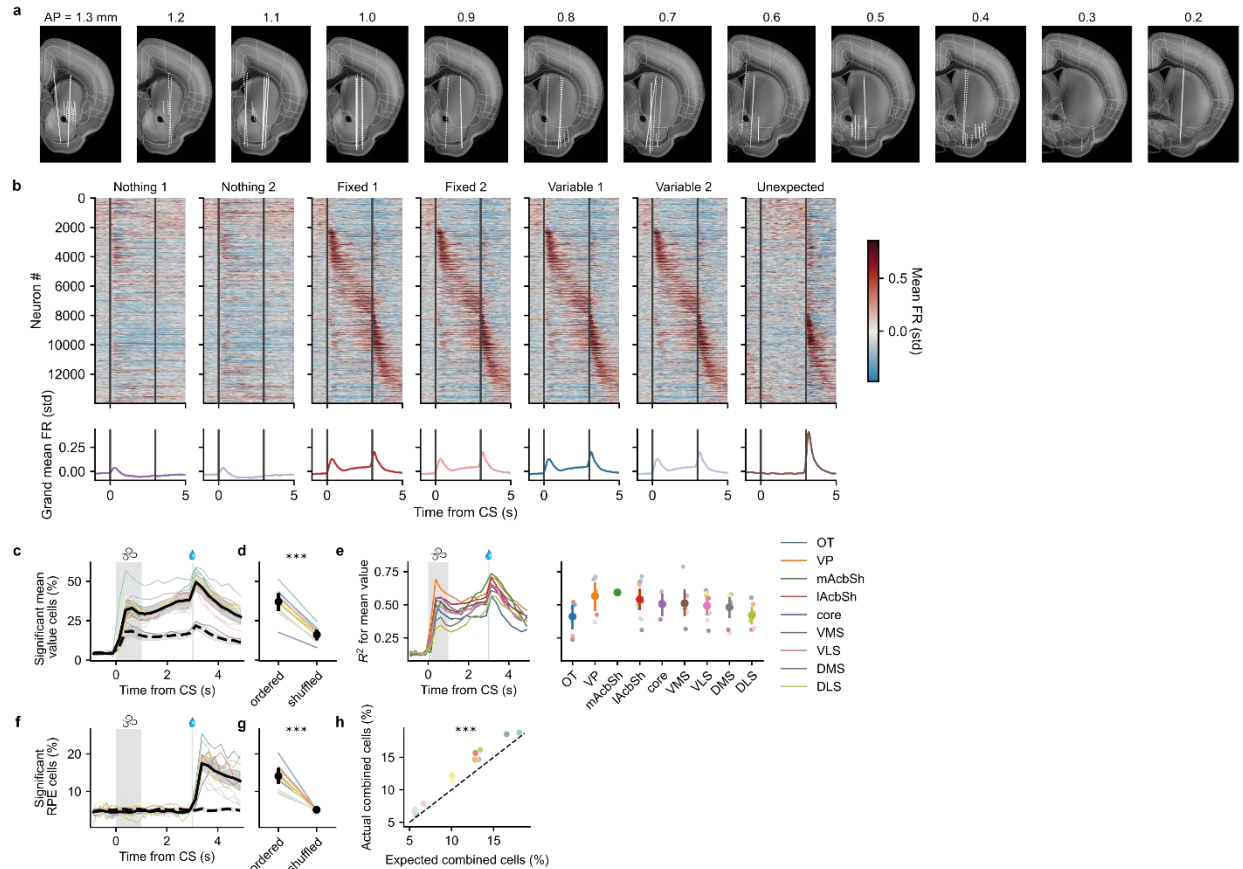


1654 **Fig. 6 | Causal contributions of D1 and D2 MSNs to REDRL.** **a**, Schematic illustration of
 1655 optogenetics experiments. For excitation, a Cre-dependent virus containing the ultrasensitive
 1656 excitatory opsin CoChR was injected into the lateral striatum at four separate depths. For
 1657 inhibition, we used transgenic animals expressing the inhibitory opsin GtACR1, also in a Cre-
 1658 dependent manner. Cre was delivered transgenically by way of *Drd1-Cre* or *Adora2a-Cre* mice,
 1659 and a tapered fiber was implanted in the lAcbSh. **b**, The trial structure in these experiments was
 1660 identical to the recording experiments, except that stimulation was delivered throughout the
 1661 duration of the Trace period. Licking was again quantified during the Late Trace period, 1–0 s
 1662 before the outcome, to avoid counting any artifactual licking around stimulation onset. The laser
 1663 was pulsed for CoChR-based excitation and continuous for GtACR1-based inhibition. **c**,
 1664 Approach for simulating the effects of optogenetic inhibition in the REDRL model. Within each
 1665 group of panels (Nothing, Fixed, and Variable), the left column shows the predicted D1 (yellow)
 1666 and D2 (green) activities for the No Manipulation (faded circles) and Manipulation (“x”s)
 1667 conditions. Inhibition is simulated by clamping the relevant population to zero. The middle

1668 column portrays the resulting effect on the encoded value predictors, V_i . In the REDRL model,
1669 optimistic ($\tau > 0.5$; blue) and pessimistic ($\tau < 0.5$; purple) predictors are identified with D1 and
1670 D2 MSNs, respectively. However, since the encoding of D2 MSNs is flipped, inhibition actually
1671 drives these V_i 's positive relative to their baseline (grey). The right column illustrates the effect
1672 this change in V_i has on the encoded mean (blue and purple horizontal dashed lines), relative to
1673 the unperturbed distribution (grey histogram, with mean shown in black). The ground-truth
1674 distributions shown reflect the versions used in the manipulation experiments, where the
1675 Variable condition consisted of equally probable 0 and 8 μL rewards. **d**, Same as **c**, but for
1676 optogenetic excitation (triangles) rather than inhibition. **e**, Summary of REDRL model
1677 predictions. Each point was computed as the difference in the implied mean between the
1678 Manipulation and No Manipulation conditions, computed separately for inhibition (*left*) and
1679 excitation (*right*). **f**, Difference in Late Trace period anticipatory licking between lAcSh
1680 Manipulation and No Manipulation trials, computed within-session and then averaged across-
1681 session and within-mice (thin lines). Thick lines and shaded regions show the mean \pm 95%
1682 confidence interval across mice. To emphasize the concordance with REDRL predictions, D1
1683 and D2 manipulations are now colored blue and purple, respectively. Colored asterisks with
1684 horizontal lines denote significant differences in the effect of manipulation between trial types
1685 within the indicated genotype (D1 inhibition: $p < 0.001$ Nothing vs. Fixed, $p < 0.001$ Nothing vs.
1686 Variable; D1 excitation: $p < 0.001$ Nothing vs. Fixed, $p < 0.001$ Nothing vs. Variable; D2
1687 excitation: $p = 0.007$, Nothing vs. Fixed). Colored asterisks over single trial types indicate
1688 significant differences relative to zero for that genotype (D2 inhibition: $p < 0.001$ Nothing, $p =$
1689 0.002 Fixed, $p < 0.001$ Variable; D1 excitation: $p < 0.001$ Nothing; D2 excitation, $p = 0.032$
1690 Nothing). Black asterisks over single trial types indicate significant differences between
1691 genotypes (inhibition: $p = 0.001$ Fixed, $p = 0.005$ Variable; excitation: $p < 0.001$ Nothing). **g**,
1692 Summary panel showing the mean coefficient of determination for each model, used to predict
1693 the average difference in licking for each trial type without any intercept term.

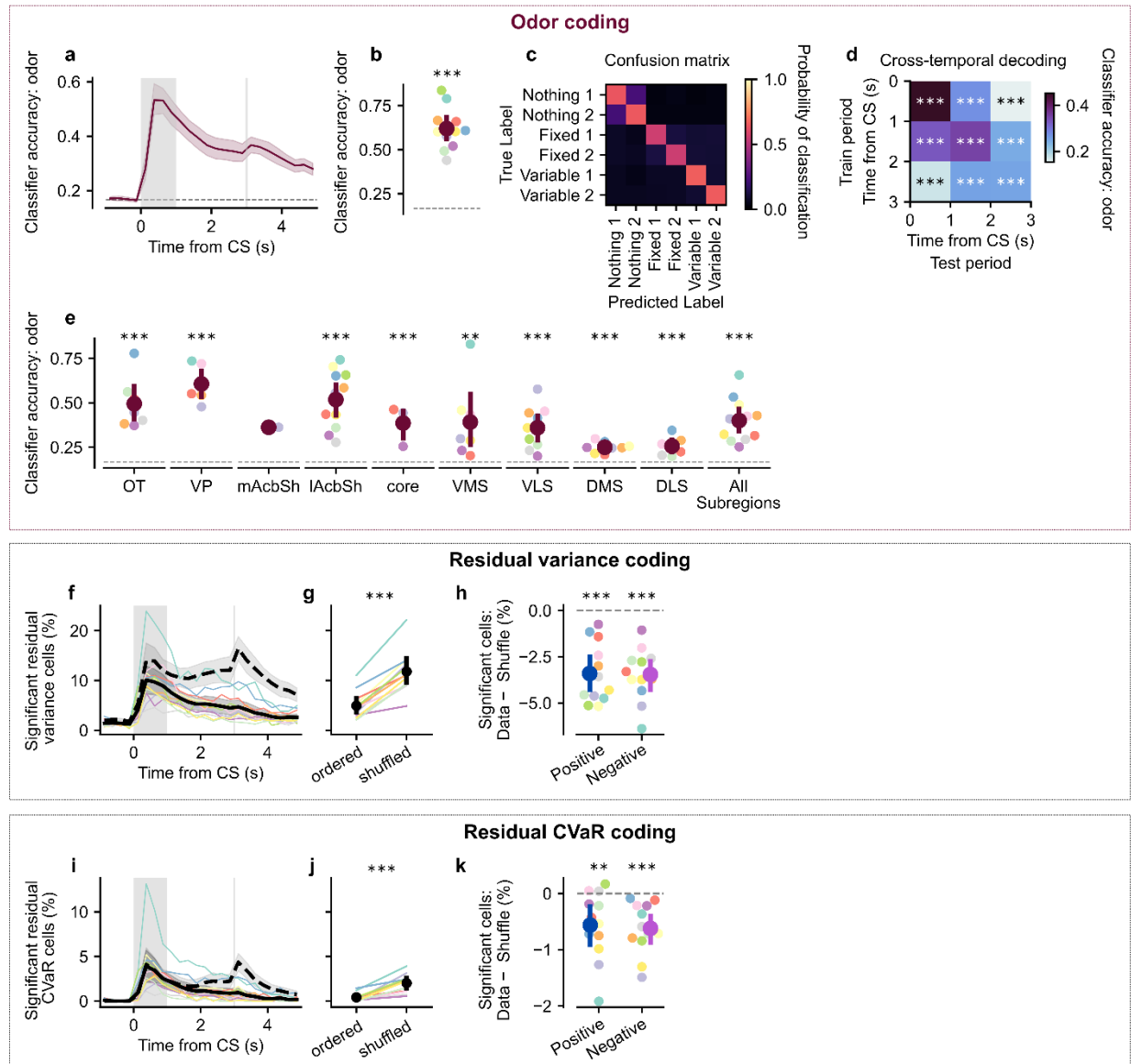


1694 **Extended Data Fig. 1 | Behavioral classification analysis.** **a**, Odors corresponding to the same
1695 distribution were treated as the same class. This is illustrated for the case of Fixed vs. Variable
1696 classification, with the background shading (yellow vs. grey) indicating the target for the
1697 classifier. **b**, Schematic of behavioral classification. On each validation fold, whisking, running,
1698 pupil area, licking, and the top 50 face motion energy PCs in the training set were z-scored and
1699 then passed to a support vector classifier (SVC) with a linear kernel, which predicts the
1700 associated distribution. **c**, Schematic of orthogonality analysis. The weights learned by the SVC
1701 define a vector orthogonal to the hyperplane that best separates distributions. A separate vector
1702 can be defined by regressing the mean reward (“Value direction”) of each trial against their
1703 corresponding behavioral regressors. While the SVC hyperplane considers only four odors at a
1704 time, the regression direction takes into account all six odors. **d**, Cosine similarity between the
1705 classifier weight vector and the Value direction. Any differences in behavior between Fixed and
1706 Variable trials are orthogonal to Value (relative to chance level of 0: $p < 0.001$ for Nothing vs.
1707 Fixed, $p < 0.001$ for Nothing vs. Variable, $p = 0.154$ for Fixed vs. Variable). **e**, Spatial masks
1708 corresponding to face motion energy PCs in an example session, sorted by variance explained.
1709 Successive PCs emphasize finer and finer aspects of mouse whisking, sniffing, and licking
1710 behavior.



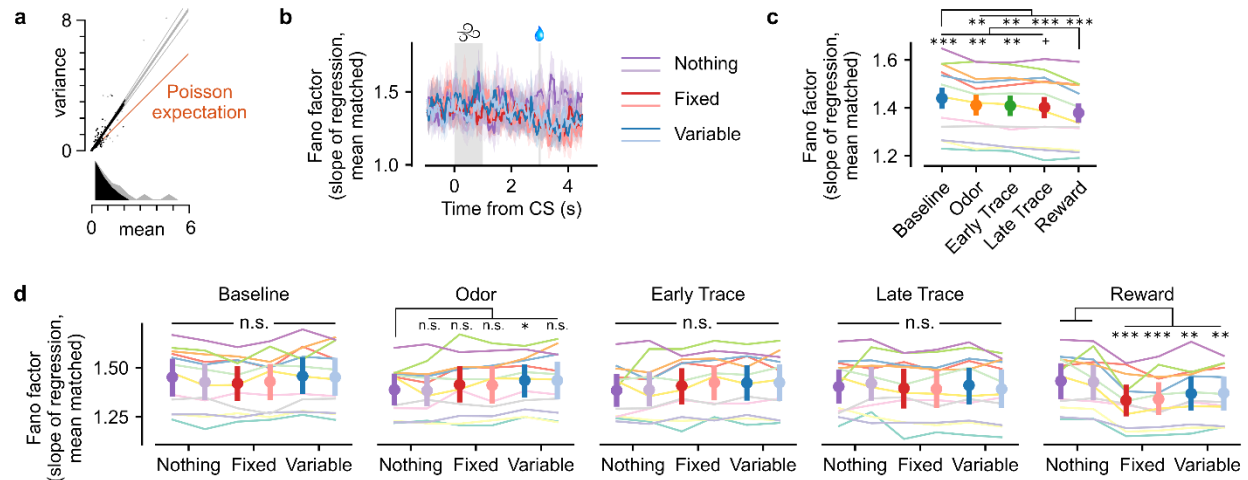
1711 **Extended Data Fig. 2 | Value and RPE coding across the striatum.** **a**, Serial coronal sections
 1712 showing recording sites of probe insertions (white dotted lines), registered to the Allen Common
 1713 Coordinate Framework. **b**, *Top*, heatmaps showing average z-scored firing rate in response to
 1714 each odor for each neuron. Neurons were sorted according to the time of peak activity when
 1715 averaged on half of Variable 2 odor trials, and then plotted in this same order for the remainder
 1716 of trials, grouped by trial type. The seventh and final trial type corresponds to Unexpected
 1717 rewards, which were not preceded by an odor. *Bottom*, grand average z-scored firing rate across
 1718 all neurons. **c**, Fraction of neurons that significantly correlate with mean reward, computed
 1719 separately in non-overlapping 250 ms time bins. Each mouse is shown in a different color, with
 1720 the mean \pm 95% confidence interval across mice shown in solid black. Dashed line is the average
 1721 across mice after shuffling the mapping between odors and distributions, thereby accounting for
 1722 pure odor coding. **d**, Average percentage of significant cells during the Late Trace period ($p <$
 1723 0.001 , paired samples t -test). **e**, *Left*, cross-validated R^2 predicting the mean reward on each trial
 1724 as a function of striatal subregion, computed separately in non-overlapping 250 ms time bins. To
 1725 ensure fair comparison across subregions, we for each animal generated multiple pseudo-
 1726 populations of 40 neurons each by repeatedly sampling without replacement neural
 1727 subpopulation across session boundaries until there were fewer than 40 neurons remaining.
 1728 Animals with fewer than 40 neurons in the given region were excluded. Lines show averages
 1729 across mice for each subregion. *Right*, average R^2 over the Late Trace period. Smaller dots show

1730 averages across pseudo-populations for each mouse with at least 40 neurons in that region. **f**,
1731 Same as **c**, except showing the fraction of neurons that significantly correlate with reward
1732 prediction error (RPE), defined as the difference between actual and expected reward. **g**, Same as
1733 **d**, except showing the average percentage of significant cells during the Outcome period, 0–1 s
1734 after reward delivery ($p < 0.001$). **h**, The actual fraction of cells in each mouse that significantly
1735 correlated with both mean value and RPE was compared to the product of the individual
1736 fractions for mean and RPE-coding cells (the predicted fraction assuming independence; $p <$
1737 0.001, paired samples t -test).



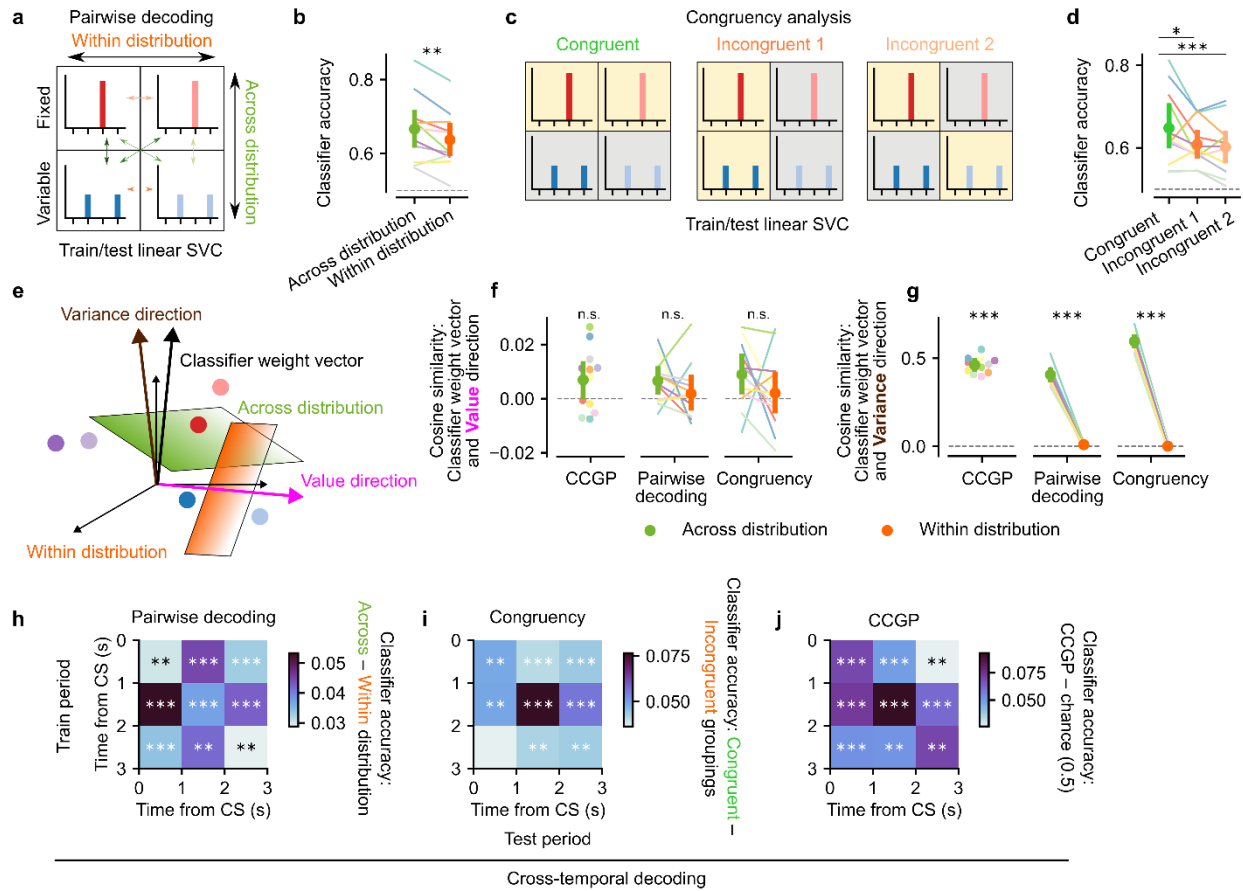
1738 **Extended Data Fig. 3 | Odor and residual variance coding in the striatum.** **a**, Decoding
 1739 accuracy across time of a multinomial logistic regression classifier decoding odor identity. **b**,
 1740 Quantification of **a** during the Odor period ($p < 0.001$ relative to chance level of 1/6). **c**,
 1741 Confusion matrix for odor decoding during the odor period shows high decoding accuracy for all
 1742 odors, with relatively higher confusability for odors with the same mean. **d**, Cross-temporal
 1743 decoding reveals that odor decoding is stable across time, allowing a classifier trained e.g. on
 1744 Late Trace period activity to generalize well above chance to the Odor period, and vice versa (all
 1745 p 's < 0.001 relative to chance level of 1/6). **e**, Pseudo-population odor decoding across
 1746 subregions (see Methods section titled "Comparisons across subregions, hemispheres, and
 1747 genotypes"). OT, olfactory tubercle; VP, ventral pallidum; mAcbSh, medial nucleus accumbens
 1748 shell; lAcbSh, lateral nucleus accumbens shell; core, nucleus accumbens core; VMS,
 1749 ventromedial striatum; VLS, ventrolateral striatum; DMS, dorsomedial striatum; DLS,

1750 dorsolateral striatum ($N = 1$ mouse for mAcSh, $p = 0.006$ for VMS, all other p 's < 0.001). **f**,
1751 Same as Extended Data Fig. 2c, except showing the fraction of neurons that significantly
1752 correlate with variance, after regressing out the contribution of mean reward coding separately
1753 for each time bin. **g**, Average percentage of significant Residual Variance cells during the Late
1754 Trace period is *less* than would be predicted from odor coding alone ($p < 0.001$, paired samples
1755 t -test). **h**, Same as Fig. 3m, except for Residual Variance coding. Fraction is lower than chance
1756 for both positive- and negative-coding cells ($p < 0.001$, paired samples t -test). **i-k**, Same as **f-h**,
1757 but for conditional value at risk (CVaR), a common risk measure used in finance and
1758 reinforcement learning^{126,176,177}, defined as the expected value within the lower α -quantile of a
1759 probability distribution. For our distributions, this will be equivalent to the mean for $\alpha > 0.5$ and
1760 equivalent to the minimum value for $\alpha < 0.5$, which differs only for the Variable distribution,
1761 where it is 2. The latter is what we plot here, after regressing out mean coding. Again, there are
1762 fewer Residual CVaR cells than would be expected from odor coding alone ($p < 0.001$, paired
1763 samples t -test) and this is true for both positive- and negative-coding cells ($p = 0.009$ and $p <$
1764 0.001 , respectively, paired samples t -test).



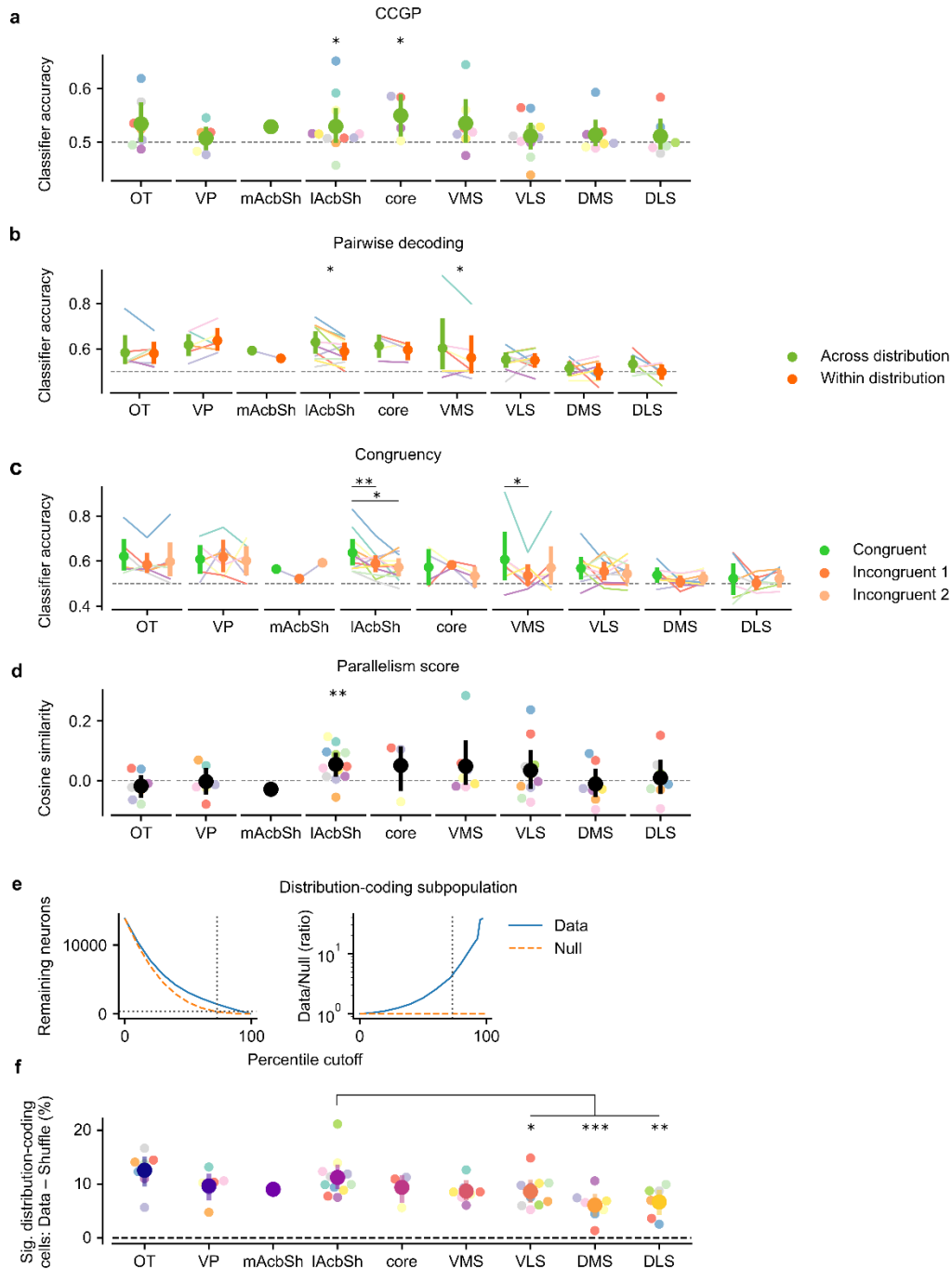
1765 **Extended Data Fig. 4 | Sampling-based codes are inconsistent with striatal activity patterns.**

1766 **a**, Illustration of how the mean-matched Fano factor was computed¹⁷⁸. Spike counts were
 1767 computed in 100 ms bins for each trial. The mean and variance (across trials) of that count then
 1768 contributed one data point to the scatter plot. Grey dots depict all neurons from an example
 1769 session, time bin (centered 200 ms after odor onset), and odor (Variable 2). The grey line is the
 1770 regression fit to all data, constrained to pass through zero and weighted according to the
 1771 estimated s.e.m. of each variance measurement. Black dots are the data points preserved by mean
 1772 matching at each time point, to eliminate the possibility that differences across time are driven by
 1773 differences in firing rates, which could in principle violate the Poisson assumption. This
 1774 transforms the distribution of mean counts from the grey to the black distribution. The regression
 1775 slope for the mean matched data is plotted as the black line. Finally, the Poisson expectation of
 1776 equal mean and variance is plotted in orange, with a slope of one. This procedure was performed
 1777 independently on each session, time bin, and trial type. **b**, Time course of the computed mean-
 1778 matched Fano factor (\pm 95% confidence interval) for the example session shown in **a**. That is, the
 1779 slope of black line in **a** is the height of the light blue, Variable 2 line in **b** 200 ms after CS onset.
 1780 **c**, Quantification of mean matched Fano factor across second-long time periods. Consistent with
 1781 cortical observations¹⁷⁸, we see a quenching of variability upon CS onset (Baseline: $p = 0.002$,
 1782 0.001 , < 0.001 , < 0.001 relative to Odor, Early Trace, Late Trace, and Reward periods), and
 1783 another one upon reward delivery (Reward: $p < 0.001$, $= 0.002$, 0.006 , 0.053 for Baseline, Odor,
 1784 Early, and Late Trace periods). **d**, Quantification of mean matched Fano factor across trial types,
 1785 shown separately for each time period. In general, there is no tendency for Variable odors to
 1786 elicit strong and sustained increases in variability, as would be predicted by sampling-based
 1787 codes (Baseline, Odor, Early and Late Trace: all p 's > 0.05 , except Nothing 1 vs. Variable 1 for
 1788 Odor: $p = 0.032$ uncorrected). However, reward delivery specifically drives yet another decrease
 1789 in variability (Nothing 1: $p = 0.570$ for Nothing 2; $p < 0.001$ for Fixed odors; $p = 0.002$ for
 1790 Variable odors).



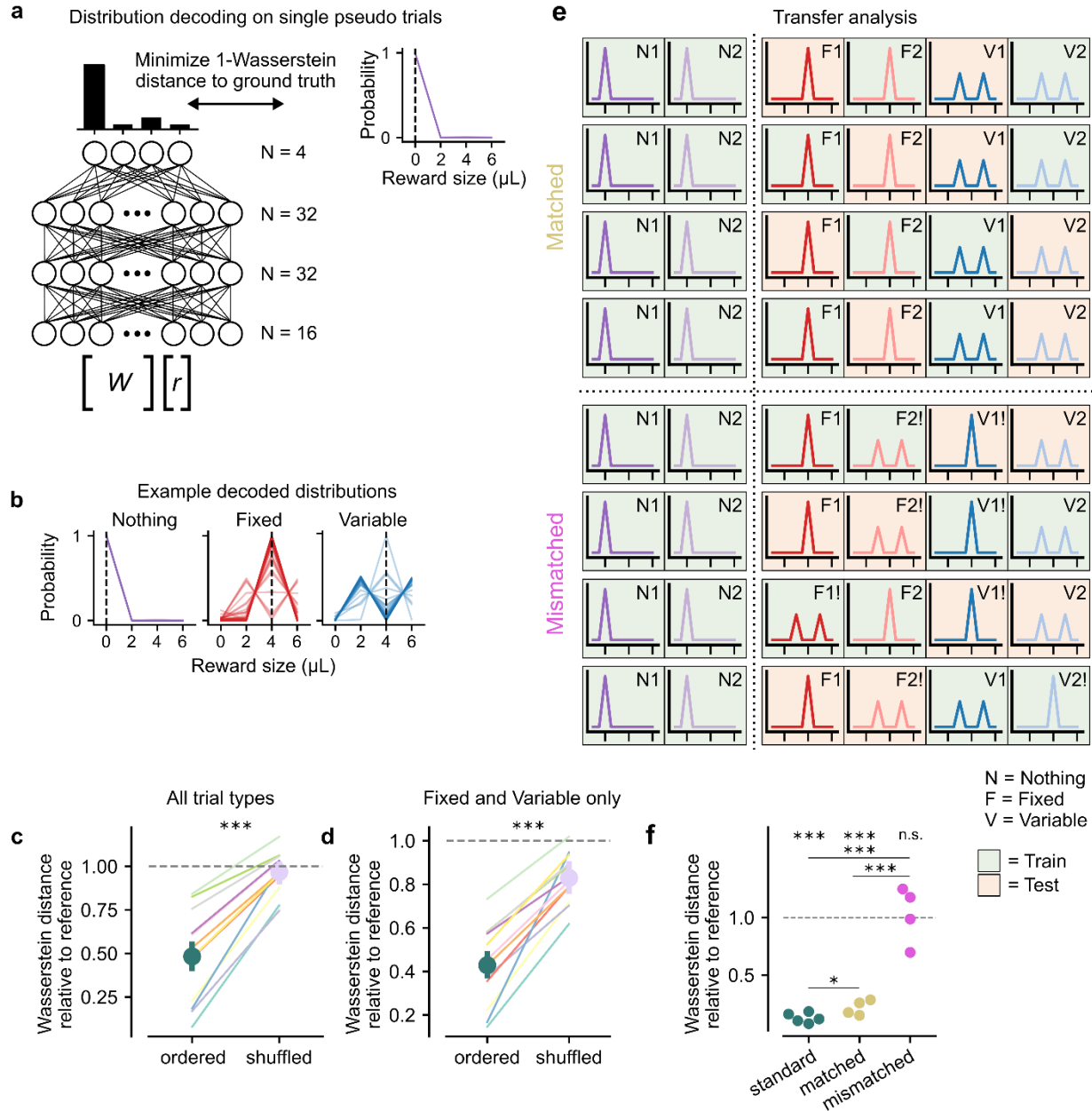
1791 **Extended Data Fig. 5 | Distributional coding is robust, orthogonal to value, and consistent**
 1792 **across time.** **a**, Schematic of pairwise decoding analysis. Linear SVCs were trained
 1793 individual Fixed and Variable odors, two at a time. This resulted in six possible dichotomies,
 1794 four of which encompassed one Fixed and one Variable odor (green arrows; “Across
 1795 distribution”) and two of which compared odors cuing the same exact distribution (orange
 1796 arrows; “Within distribution”). **b**, Pairwise decoding during the Late Trace period was
 1797 significantly better for across- than within-distribution pairs, consistent with distributional but
 1798 not traditional RL ($p = 0.001$). **c**, Schematic of congruency analysis, which considered all four
 1799 Fixed and Variable odors simultaneously. In the Congruent grouping, both Fixed odors were
 1800 assigned to one class (yellow background) and both Variable odors were assigned to the other
 1801 class (grey background), just as was done for behavioral decoding. By contrast, in the
 1802 Incongruent groupings, class assignments cut across Fixed and Variable distributions. **d**,
 1803 Classifier accuracy in the Late Trace period was higher for Congruent than Incongruent pairs,
 1804 again consistent with distributional but not traditional RL (Congruent: $p = 0.028$ vs. Incongruent
 1805 1, $p < 0.001$ vs. Incongruent 2). **e**, Schematic illustrating the classifier weight vector (normal to
 1806 the separating hyperplane for across- or within-distribution classifications) and the regression
 1807 weight vector (for Value or Variance). **f**, Quantification of cosine similarity between the
 1808 classifier weight vector and the Value direction shows that the vectors are not significantly
 1809 different from orthogonal (CCGP: $p = 0.071$ relative to chance value of 0; Pairwise: $p = 0.797$

1810 Across- vs. Within-distribution absolute cosine similarity; Congruency: $p = 0.493$ Across- vs.
1811 Within-distribution absolute cosine similarity). **g**, Same as **f**, but for Variance rather than Value
1812 direction ($p < 0.001$ for all comparisons). **h-j**, Cross-temporal decoding for the pairwise,
1813 congruency, and CCGP analyses. Distributional RL is favored during every time period between
1814 odor onset and reward delivery, and decoders trained during one period almost always generalize
1815 to other time periods.



1816 **Extended Data Fig. 6 | Distributional coding is strongest in the lAcSh.** **a**, Pseudo-population
 1817 CCGP across subregions (relative to chance level of 0.5: $p = 0.059, 0.473, 0.044, 0.017, 0.088,$
 1818 $0.346, 0.257, 0.407,$ and 0.133 for OT, VP, mAcSh, lAcSh, core, VMS, VLS, DMS, and DLS,
 1819 respectively. Same order applies to all statistics in this figure). Pseudo-populations were
 1820 constructed as in Extended Data Fig. 3e. **b**, Pseudo-population pairwise decoding across
 1821 subregions (Across- vs. Within-distribution: $p = 0.861, 0.344, 0.883, 0.010, 0.409, 0.040, 0.882,$
 1822 $0.482, 0.106$). **c**, Pseudo-population congruency analysis across subregions (Congruent vs.
 1823 Incongruent 1: $p = 0.097, 0.817, 0.744, 0.007, 0.832, 0.047, 0.523, 0.138, 0.523$; Congruent vs.

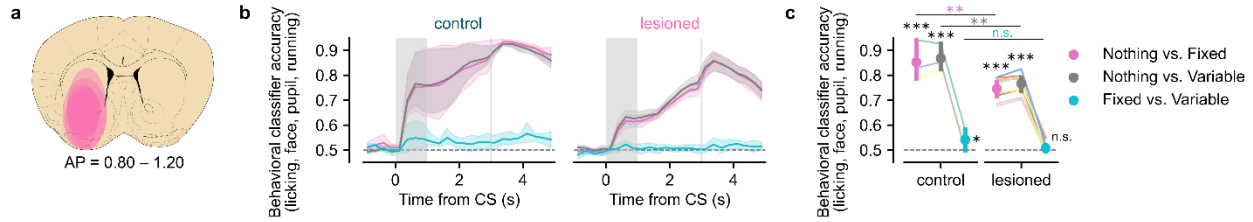
1824 Incongruent 2: $p = 0.306, 0.760, 0.815, 0.010, 0.473, 0.177, 0.316, 0.486, 0.985$). **d**, Parallelism
1825 score across subregions (relative to chance level of 0: $p = 0.300, 0.878, 1.00, 0.001, 0.229, 0.243,$
1826 $0.273, 0.615, 0.764$). **e**, *Left*, fraction of neurons with classifier coefficients above the percentile
1827 cutoff for all three (CCGP, pairwise, and congruency) analyses. Horizontal dotted line indicates
1828 level at which 2.5% of null coefficients fell above the cutoff; this was the 73rd percentile
1829 (vertical dotted line), and retained 11.43% of neurons. *Right*, ratio of data to null coefficients
1830 falling above the cutoff (log scale). **f**, Fraction of distribution-coding cells in each subregion.
1831 This fraction is significantly higher in the lAcbSh than in more dorsal subregions (relative to
1832 lAcbSh: $p = 0.339, 0.285, 0.473, 0.274, 0.071, 0.038, 0.001$ for OT, VP, mAcbSh, core, VMS,
1833 VLS, and DLS, respectively; $p < 0.001$ for DMS).



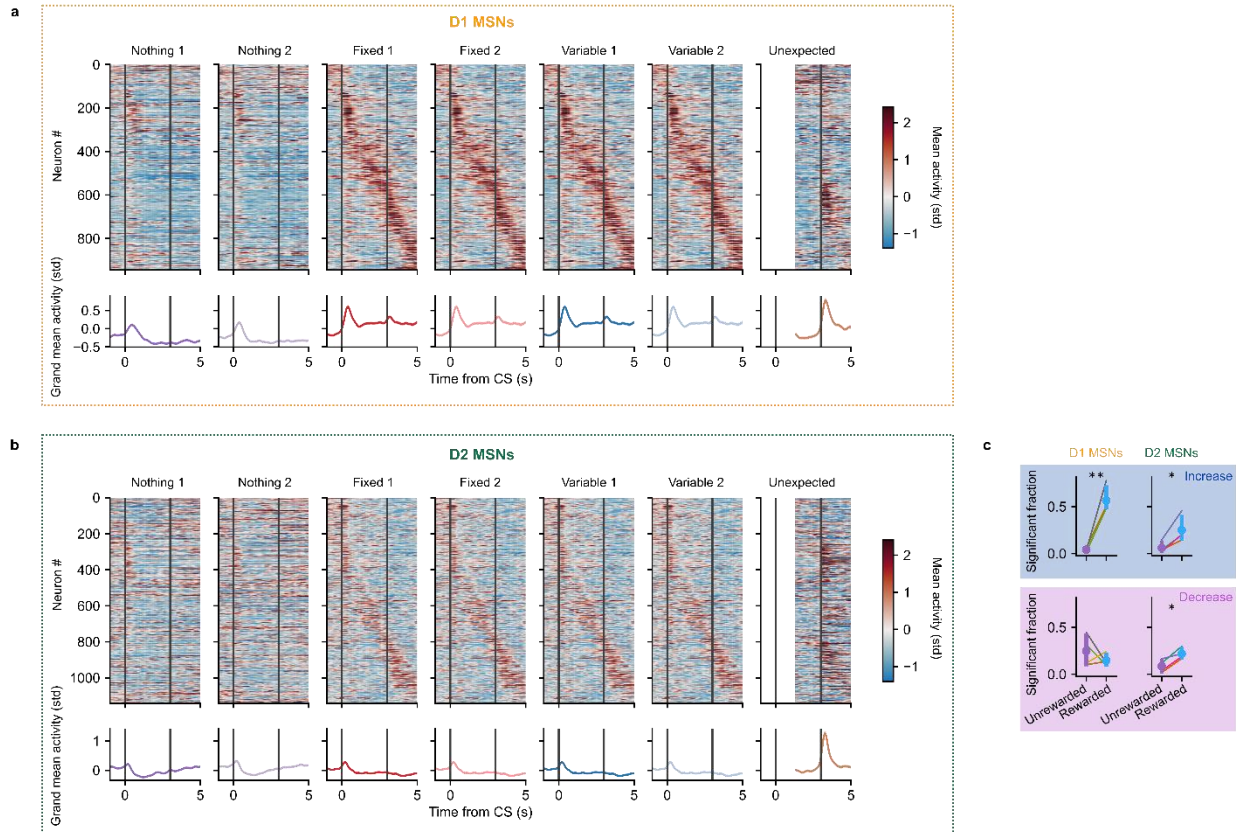
1834 **Extended Data Fig. 7 | Artificial neural network-based distribution decoding captures**
 1835 **information beyond the mean.** **a**, ANN schematic. Single-trial spike counts from the
 1836 distribution-coding subpopulation r were linearly mapped into 16 dimensions by the trainable
 1837 matrix W and then fed through the network (see Methods). After a final layer, a softmax function
 1838 transformed activations into a properly-normalized probability distribution, whose 1-Wasserstein
 1839 distance to ground truth was minimized with stochastic gradient descent. **b**, Example decoded
 1840 distributions from the test set, shown as line plots to distinguish individual pseudo-trials. **c**,
 1841 Wasserstein distance relative to reference for the ANN trained on all six trial types, with and
 1842 without shuffling odor-distribution mappings ($p < 0.001$ ordered vs. shuffled; $p < 0.001$ ordered
 1843 relative to chance value of 1; $p = 0.350$ shuffled relative to chance value of 1). **d**, Same as **c**, but

1844 for ANN trained on only the Rewarded odors, which shared the same mean ($p < 0.001$ ordered
1845 vs. shuffled, ordered relative to chance value of 1, and shuffled relative to chance value of 1). **e**,
1846 Schematic depicting setup for transfer analysis. Four trial types, including both Nothing odors,
1847 were used for training (green background), and the other two were used for testing (orange
1848 background). Matched pairings veridically assigned odors to distributions, while mismatched
1849 pairings used either only Fixed or only Variable odors for training while assigning one member
1850 per training pair and one member per testing pair to the opposite distribution (indicated by the
1851 exclamation mark). There were four possible ways to draw the matched dichotomies, all of
1852 which are shown (rows). For the mismatched dichotomies, the test labels could be flipped
1853 arbitrarily, so only one possibility (the F2 and V1 distributions swapped for testing) is shown for
1854 each training set. **f**, Wasserstein distance relative to reference for standard, matched, and
1855 mismatched settings. Standard is identical to analysis shown in **c**, except that for this decoder,
1856 neurons from all mice were pooled. Matched transfer yields distributions that are nearly as
1857 accurate as training with all six trial types ($p < 0.001$ for matched vs. mismatched and standard
1858 vs. mismatched, independent samples t -test; $p = 0.043$ for standard vs. matched, independent
1859 samples t -test; $p < 0.001$ for standard and matched relative to chance value of 1, one-sample t -
1860 test; $p = 0.836$ for mismatched relative to chance value of 1, one-sample t -test).

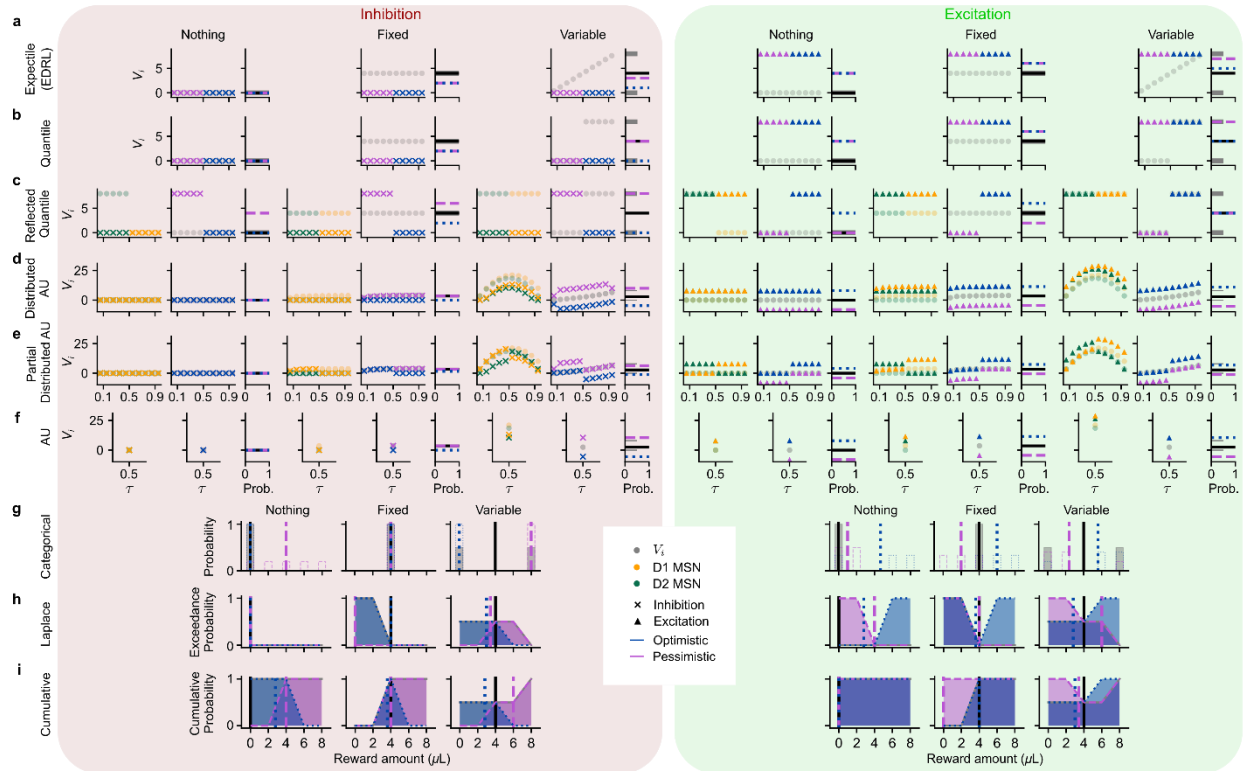
1870 level of τ (grey dots) approximates each expectile, and their sum relates to the spread of the
1871 distribution. This drives maximal activity in response to Variable odors, which is why they
1872 separate out most clearly along PC 1. **e**, Same as **d**, but for a reduced version in which only a
1873 single pair of value predictors are learned with balanced positive and negative learning rates⁶⁶ (τ
1874 = 0.5). **f**, Same as **a**, but for a categorical code in which distributions are encoded as a
1875 histogram⁷¹. Each neuron is imagined to correspond to a single reward bin, with its firing rate
1876 proportional to the height of that bin. **g**, Same as **f**, but for a Laplace code⁸³. In the limit of
1877 infinitely steep reward sensitivities for the teaching signal, these value predictors converge to the
1878 probability that the reward delivered exceeds some threshold reward amount, the “exceedance
1879 probability.” This is simply 1 minus the CDF of the probability distribution in question. Neural
1880 activities are taken to be proportional to this 1 – CDF value. **h**, Same as **g**, but for a population of
1881 neurons that flips the encoding, and so is directly proportional to the CDF. **i**, A hypothetical
1882 “distributional” code in which each neuron’s firing rate linearly correlates with either reward
1883 mean (*left*) or variance (*right*). **j**, Each trial type, replotted in mean–variance space. From this
1884 picture, it is clear that for this particular set of reward distributions, Fixed odors will be located at
1885 the midpoint between Nothing and Variable odors along PC 1, though altering the ratio of mean-
1886 to variance-coding neurons will move Fixed odors left or right along PC 1. Different sets of
1887 reward distributions could lead to different geometries. **k–m**, Qualitative features of each code in
1888 **a–i** plus random noise. REDRL predictions from Fig. 3 are included in the box on the second-to-
1889 last line, for comparison. **k**, PCA projection for each code. Only quantile-like codes give rise to
1890 the pattern observed in the data. **l**, Percentage of simulated predictors that significantly correlate
1891 with mean reward either positively (blue) or negatively (purple) for each code type. Only the
1892 reflected and categorical codes have a substantial fraction of both types of cells. In practice the
1893 positive-coding predictors are optimistic and the negative-coding predictors are pessimistic. **m**,
1894 Hypothetical activity in response to each distribution, averaged separately over optimistic (blue)
1895 and pessimistic (purple) predictors for each code type. Only the reflected codes and AU model
1896 predict a noticeable uptick in Variable relative to Fixed odors.



1897 **Extended Data Fig. 9 | Quantification of 6-OHDA lesion extent, location, and behavior. a,**
1898 Consensus heat map of all five animals' lesion locations. 6-OHDA was injected in the lAcSh
1899 but diffused into the VLS, so we considered both regions to be lesioned. We excluded OT,
1900 despite the fact that it was often lesioned, because it is not physically contiguous and showed
1901 weaker evidence of distributional coding in control animals. **b,** Behavioral decoding analysis
1902 comparing fully intact animals ($N = 3$) and unilaterally lesioned ($N = 9$) animals across time. For
1903 this analysis, animals were considered lesioned if they had received any 6-OHDA injection, even
1904 if that hemisphere was never recorded or was mistargeted relative to Neuropixels recording
1905 location. **c,** Quantification of behavioral classifier accuracy during the Late Trace period. While
1906 across-mean behavioral decoding was stronger in the control than the lesioned animals (effect of
1907 lesion: $p = 0.006, 0.001, 0.173$ for Nothing vs. Fixed, Nothing vs. Variable, and Fixed vs.
1908 Variable, respectively), both groups of animals clearly learned the task and had above-chance
1909 across-mean decoding ($p < 0.001$ compared to chance level of 50% for both Nothing vs. Fixed
1910 and Nothing vs. Variable in control as well as lesioned animals). Interestingly, Fixed vs. Variable
1911 classification was also weakly significant ($p = 0.032$ relative to chance level of 50%) for fully
1912 intact control animals, providing behavioral evidence that they did in fact learn this distinction.

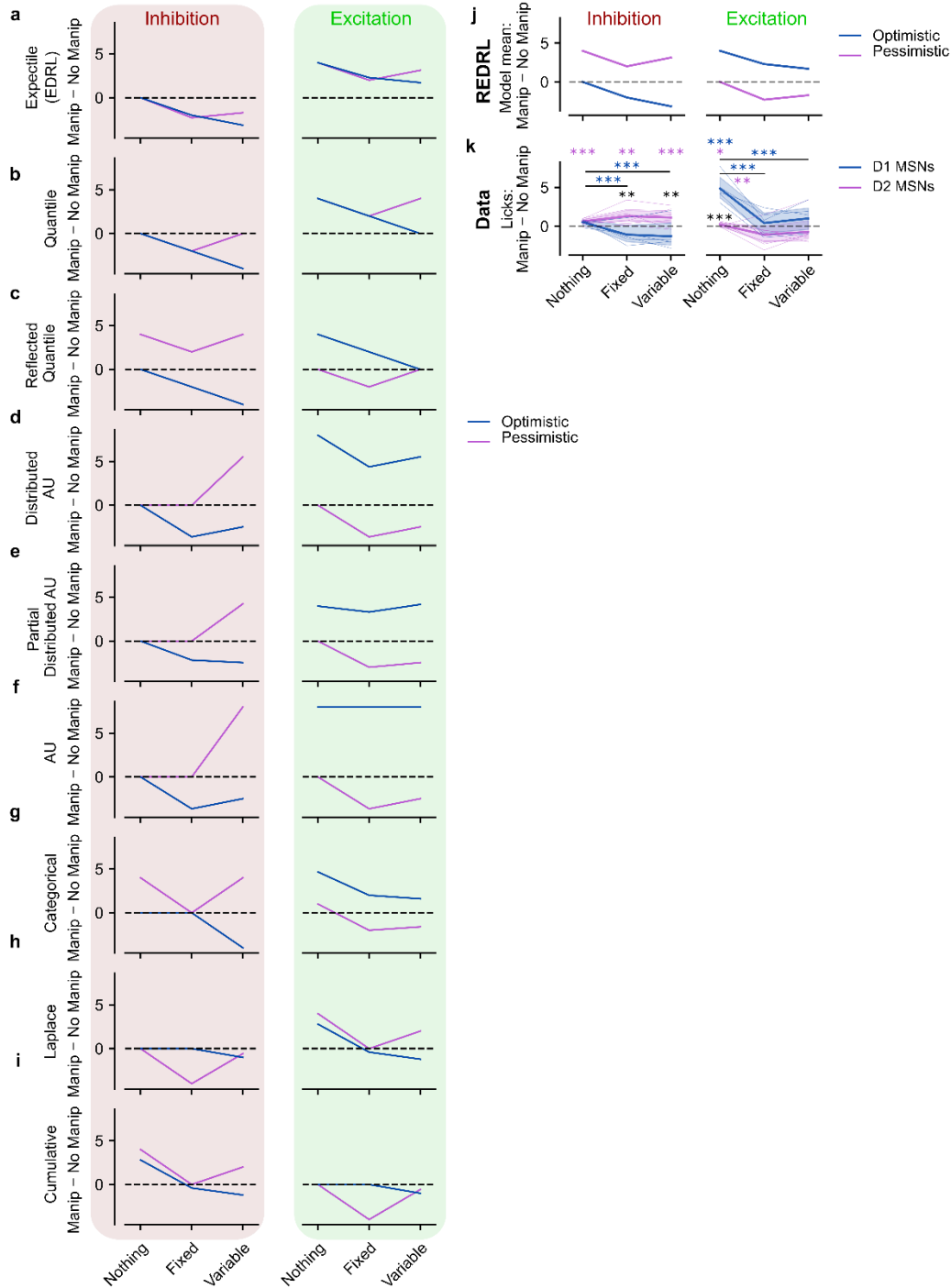


1913 **Extended Data Fig. 10 | Additional data for two-photon calcium imaging. a**, D1 MSN
 1914 activity. *Top*, heatmaps showing average z-scored deconvolved calcium activity in response to
 1915 each odor for each neuron, as in Extended Data Fig. 2b. Unexpected reward trials were cropped
 1916 on the left to include only continuous acquisitions. *Bottom*, grand average z-scored deconvolved
 1917 calcium activity across all neurons. **b**, Same as **a**, but for D2 MSN activity. **c**, Fraction of
 1918 neurons whose Late Trace activity increased (*top*) or decreased (*bottom*) relative to Baseline,
 1919 shown separately for D1 (*left*) and D2 (*right*) MSNs and Unrewarded (Nothing) versus
 1920 Rewarded (Fixed and Variable) odors (*x*-axis); these trial types were pooled before analysis. As
 1921 expected, a larger fraction of D1 MSNs increases to Rewarded rather than Unrewarded odors (p
 1922 = 0.006), while there is no difference in the fractions that decrease ($p = 0.423$). Meanwhile, for
 1923 D2 MSNs, a significantly greater fraction of neurons change their activity on Rewarded
 1924 compared to Unrewarded trials, by either increasing ($p = 0.022$) or decreasing ($p = 0.016$) their
 1925 activity relative to Baseline. Asterisks and p -values report the results of paired t -tests on
 1926 Rewarded vs. Unrewarded fractions across mice.



1927 **Extended Data Fig. 11 | Additional detail for distributional model manipulations. a,**
 1928 Schematic showing how optogenetic perturbations were simulated for an expectile code (from
 1929 EDRL). Optimistic (blue) or pessimistic (purple) predictors were shifted from their original
 1930 values (semi-transparent grey circles) and clamped to low or high values to mimic inhibition
 1931 (*left*, “x”s) or excitation (*right*, triangles), respectively. Panels on the right depict the ground-
 1932 truth reward distribution, its mean (black line), and the means of the manipulated sets of value
 1933 predictors (blue or purple dashed lines). **b**, Same as **a**, but for a quantile rather than expectile
 1934 code. **c**, Same as **b**, but for a reflected quantile code. The additional, leftmost panel for each
 1935 distribution depicts the activity of D1 (yellow) and D2 (green) MSNs at baseline (semi-
 1936 transparent circles) and after manipulations (opaque “x”s and triangles). These are what are
 1937 directly clamped by the simulated optogenetic inhibition or excitation. As a result, the effect on
 1938 the implied value predictors (middle panel) corresponding to D2 MSNs are of opposite sign, as is
 1939 the change in predicted mean (right panel). **d**, Same as **c**, but for the Distributed Actor
 1940 Uncertainty (AU) model. Since D1 and D2 MSN activities in this model can exceed the
 1941 maximum reward value, the left panel shows that perturbations were simulated by adding or
 1942 subtracting a fixed amount from each activity level (opaque “x”s and triangles) relative to
 1943 baseline (semi-transparent circles). The middle panel plots the resulting value predictors,
 1944 computed as the pointwise differences between D1 and D2 MSN activities, for pessimistic
 1945 (purple) and optimistic (blue) manipulations in comparison to baseline (grey semi-transparent
 1946 circles). **e**, Same as **d**, except that only the optimistic or pessimistic half of MSNs were
 1947 manipulated to simulate perturbations of D1 or D2 MSNs, respectively. **f**, Same as **d**, except for
 1948 the original Actor Uncertainty (AU) model in which there is only one pair of value predictors

1949 with balanced learning rates ($\tau = 0.5$). **g**, Schematic showing how optogenetic perturbations were
1950 simulated for a categorical code (from CDRL), which effectively represents the reward
1951 distribution using a histogram. Pessimistic (0, 2 μ L; purple) or optimistic (6, 8 μ L; blue) bins
1952 were clamped to 0 or 1 to simulate inhibition or excitation, respectively, relative to baseline
1953 (grey). The resulting distributions were normalized to sum to one (see Methods). Dashed vertical
1954 lines show the means of the ground-truth (black) and manipulated distributions. **h**, Same as **g**,
1955 except for a Laplace code⁸³ in which each neuron corresponds to the height of $1 - \text{CDF}$ at a
1956 particular point. While the baseline case is always monotonically decreasing, simulated
1957 excitation or inhibition can change this. Means were computed by differentiating and then
1958 normalizing (see Methods). **i**, Same as **h**, except for a cumulative code where each neuron
1959 corresponds to the height of the CDF at a particular point.



1960 **Extended Data Fig. 12 | Summary of alternative model predictions. a-i**, Predicted difference
 1961 in mean reward due to inhibition (*left*) and excitation (*right*) for each of the alternative models in
 1962 Extended Data Fig. 11. **j**, REDRL model predictions for mean reward, copied from Fig. 6e, for
 1963 comparison. **k**, Actual differences in licking, copied from Fig. 6f, for comparison.