

# 1 **The role of prospective contingency in the control of behavior and** 2 **dopamine signals during associative learning**

3

4 Lechen Qian<sup>1,2,4</sup>, Mark Burrell<sup>1,2,4</sup>, Jay A. Hennig<sup>2,3</sup>, Sara Matias<sup>1,2</sup>, Venkatesh. N. Murthy<sup>1,2</sup>,  
5 Samuel J. Gershman<sup>2,3</sup>, Naoshige Uchida<sup>1,2,5</sup>

6

7 <sup>1</sup> Department of Molecular and Cellular Biology, Harvard University, Cambridge, MA, USA

8 <sup>2</sup> Center for Brain Science, Harvard University, Cambridge, MA, USA

9 <sup>3</sup> Department of Psychology, Harvard University, Cambridge, MA, USA

10 <sup>4</sup> These authors contributed equally.

11

12 <sup>5</sup> Correspondence to Naoshige Uchida ([uchida@mcb.harvard.edu](mailto:uchida@mcb.harvard.edu))

13

14

## 15 **Abstract**

16 Associative learning depends on contingency, the degree to which a stimulus predicts an outcome.  
17 Despite its importance, the neural mechanisms linking contingency to behavior remain elusive. Here we  
18 examined the dopamine activity in the ventral striatum – a signal implicated in associative learning – in a  
19 Pavlovian contingency degradation task in mice. We show that both anticipatory licking and dopamine  
20 responses to a conditioned stimulus decreased when additional rewards were delivered uncued, but  
21 remained unchanged if additional rewards were cued. These results conflict with contingency-based  
22 accounts using a traditional definition of contingency or a novel causal learning model (ANCCR), but can  
23 be explained by temporal difference (TD) learning models equipped with an appropriate inter-trial-  
24 interval (ITI) state representation. Recurrent neural networks trained within a TD framework develop  
25 state representations like our best ‘handcrafted’ model. Our findings suggest that the TD error can be a  
26 measure that describes both contingency and dopaminergic activity. [149 words]

## 27 Introduction

28

29 The ability to discern predictive relationships between different events is crucial for adaptive behaviors.  
30 Early investigations into animal learning revealed that mere contiguity between two events (“pairing”) is  
31 insufficient for establishing enduring associations. To understand this, consider Pavlovian conditioning,  
32 where an initially neutral cue (conditioned stimulus, CS) is paired with an outcome (unconditioned  
33 stimulus, US), such as an electrical shock. Through repeated pairings, animals learn to anticipate the  
34 outcome in response to the presentation of just the CS, leading to heightened conditioned responses (e.g.,  
35 freezing). Now, consider a scenario where the same number of pairings takes place, yet additional shocks  
36 occur in the absence of the CS, such that shocks happen with equal likelihood whether or not the CS is  
37 present. In such conditions, animals fail to display conditioned responses<sup>1-3</sup>. Moreover, when a CS  
38 predicts a decrease in the likelihood of the US, conditioned responses are reduced. Based on these  
39 experiments, Rescorla postulated that conditioning depends not on the contiguity between the CS and the  
40 US but rather on *contingency* – the degree to which the CS indicates an increase or decrease in the  
41 likelihood of the US occurring.

42 Contingency indicates conditional relationships between different events and is thought to be an  
43 important quantity not only in conditioning, but also in causal inference in statistics and artificial  
44 intelligence. What is a good measure of contingency, however, remains to be clarified<sup>4-7</sup>. One commonly  
45 adopted definition in psychology and causal inference is  $\Delta P$ , the difference in the probability of one event  
46 occurring in the presence or absence of another<sup>8-10</sup>. In Pavlovian settings with trial-like structures, such as  
47 the present study,  $\Delta P$  can be expressed as  $\Delta P = P(US|CS+) - P(US|CS-)$ , where 'CS+' and 'CS-'  
48 signify the presence and absence of a CS, respectively. While mere association does not inherently imply  
49 causality, these associations can give rise to perceived causal relationships, and it has been shown that the  
50 contingency ( $\Delta P$ ) correlates with its strength<sup>6,11,12</sup>. Although  $\Delta P$  provides a simple definition, its  
51 application necessitates a trial-like structure or defined time intervals within which the probabilities of  
52 events such as CS and US can be determined<sup>7,13</sup>. Likewise, some behavioral observations cannot be  
53 explained by  $\Delta P$ , leading some to argue against the usefulness of contingency in explaining behavior<sup>14</sup>.  
54 As a result, efforts have been made to better define contingency<sup>4-7</sup>.

55 Following Rescorla’s experiments discussed above, further experiments highlighted the crucial role of  
56 surprise in the establishment of associations<sup>15</sup>. To account for this, Rescorla and Wagner (1972)  
57 postulated that conditioning is driven by the discrepancy between the actual and predicted outcome  
58 (prediction error)<sup>16</sup>. Importantly, this contiguity-based model can explain the contingency degradation  
59 experiments described above, assuming that the context acts as another CS, which competes with the  
60 primary CS<sup>16</sup>. While this “cue-competition” account is attractive, and potentially replaces the classic  
61 contingency-based account, the validity of the cue-competition model remains contested<sup>17-20</sup>.

62 Like  $\Delta P$ , the Rescorla-Wagner model also assumes a trial-based structure, as it does not consider the  
63 timing of events either within or outside a trial. To address this limitation, Sutton and Barto developed the  
64 temporal difference (TD) learning algorithm, now a fundamental algorithm in reinforcement learning<sup>21,22</sup>,  
65 as a prediction error-based model of associative learning<sup>23,24</sup>. TD learning as a model of associative  
66 learning in animals finds support in the striking resemblance observed between the activity of midbrain  
67 dopamine neurons and the prediction error (TD error) used in TD learning algorithms<sup>25-29</sup>.

68 Despite the success of TD learning models in accounting for both associative learning and dopamine  
69 signals<sup>25,30</sup>, TD models has received various challenges from alternative models. For instance, a recent  
70 study<sup>31</sup> proposed an alternative model for associative learning and dopamine, called an adjusted net  
71 contingency for causal relations (ANCCR) model. As the name implies, the ANCCR model posits  
72 contingency as a key driver of associative learning and causal inference. Conventional definitions of  
73 contingency as well as TD learning models rely on “prospective” predictive relationships between cues  
74 and outcomes, i.e.  $P(US|CS)$ . By contrast in the ANCCR model, learning is driven by “retrospective”  
75 relationships, that is the probability of a stimulus (CS) given the outcome (US), or  $P(CS|US)$ . The authors  
76 argued that ANCCR implements causal inference, and that dopamine signals convey a signal for causal  
77 learning (the “adjusted net contingency”), not TD errors. Evidence supporting these ideas came from their  
78 experiments in examining dopamine signals in mice<sup>31</sup> and rats<sup>32</sup> during Pavlovian tasks in which  
79 contingency was manipulated. The validity of ANCCR, as well as interpretations of the data presented in  
80 these studies, await further examination.

81 The concept of contingency lies at the heart of learning predictive relationships. Recent work<sup>33,31</sup> has  
82 raised the novel question of whether associations are learned looking forward (prospectively) or looking  
83 backward (retrospectively), and how dopamine is involved in these processes<sup>7,31</sup>. Yet how contingency  
84 affects dopamine signals and behavior, as well as how dopamine signals relate to causal inference,  
85 remains to be determined. To address these questions, we examined behavior and dopamine signals in the  
86 ventral striatum (VS) in mice performing Pavlovian conditioning tasks while manipulating stimulus-  
87 outcome contingencies. We show that, contrary to previous claims<sup>31,32</sup>, dopamine signals could be  
88 comprehensively explained by TD learning models. Furthermore, we found that dopamine signals  
89 primarily reflected prospective stimulus-outcome relationships, and strongly violated predictions of the  
90 ANCCR model. We then discuss a conceptual framework for how dopamine signals can be related to  
91 contingency and causal inference.

92

93

## 94 **Results**

95

### 96 **Contingency degradation attenuates Pavlovian conditioned responding**

97 To study the effects of contingency in a Pavlovian setting, we developed a task for head-fixed mice in  
98 which odor cues predicted a stochastic reward (Fig. 1a, b, c). All mice ( $n = 29$ ), after being water  
99 restricted, were first trained on one reward-predicting odor (Odor A) that predicted a reward (9  $\mu$ L water)  
100 with 75% probability and one odor (Odor B) that indicated no reward. In this phase (Phase 1), Odor A  
101 trials accounted for 40% of trials, Odor B for 20%, with the remaining 40% being blank trials, in which  
102 neither odor nor reward was delivered. The timing of task events (Fig. 1b) was chosen such that the trial  
103 length was relatively constant, so we could apply the classic  $\Delta P$  definition to our design.

104 In Phase 1, Odor A has positive stimulus outcome contingency, being predictive of reward (R; Fig. 1c).  
105 This can be quantified using the commonly applied  $\Delta P$  definition of contingency: e.g.,  $\Delta P(A) =$   
106  $P(R|A+) - P(R|A-) = 0.75 - 0 = 0.75$  in Phase 1. Conversely, Odor B has a negative stimulus-  
107 outcome contingency:  $\Delta P(B) = P(R|B+) - P(R|B-) = 0 - 0.375 = -0.375$ . Consistent with these

108 contingencies, all animals developed anticipatory licking to Odor A, but not Odor B, within five training  
109 sessions (Fig. 1e).

110 In Phase 2, animals were split into groups (Fig. 1a). The first group ('Cond',  $n = 6$ ) continued being  
111 trained on the identical conditioning task from Phase 1. With no change in contingency, the behavior did  
112 not significantly change in a further five sessions of training (Fig. 1d, e).

113 The second group ('Deg',  $n = 11$ ) experienced contingency degradation. To reduce the contingency of  
114 Odor A, either  $P(R|A+)$  can be decreased or  $P(R|A-)$  can be increased. We increased  $P(R|A-)$  by  
115 introducing uncued rewards, an experimental design termed 'contingency degradation'<sup>34</sup>. Blank trials  
116 from Phase 1 were replaced with 'background water' trials in which a reward was delivered on 75% of  
117 these trials. In this condition,  $P(R|A+)$  remains unchanged at 0.75, while  $P(R|A-)$  increases to 0.5 (2  
118 out of every 3 non-Odor A trials are background water trials, of which 75% are rewarded thus  
119  $P(R|A-) = 2/3 \times 0.75 = 0.5$ ). As a result,  $\Delta P(A)$  is reduced to 0.25. Concomitant with this decreased  
120 contingency, the anticipatory licking to Odor A decreased across five sessions of Phase 2 in the Deg  
121 group ( $t_{11} = -4.78$ ,  $P = 0.00074$ , paired  $t$ -test). Moreover, Deg group animals increased licking during the  
122 inter-trial intervals (ITIs,  $t_{11} = 3.34$ ,  $P = 0.0074$ , paired  $t$ -test), potentially reflecting an increased baseline  
123 reward expectation. Additionally, the Deg group exhibited longer latencies to initiate licking and an  
124 increase in trials where mice did not lick before water delivery in Odor A trials (Extended Data Fig. 1d,  
125 e).

126 The decrease in anticipatory licking, rather than reflecting the decreased contingency, could reflect satiety  
127 effects as animals in the Deg group receive twice as many rewards per session as the Cond group. We do  
128 not believe satiety explains this effect for at least two reasons: (1) all animals still received and drank  
129 about 1 ml supplementary water after each session, and (2) in all but the first degradation session,  
130 anticipatory licking was diminished compared to Cond controls in early trials (Extended Data Fig. 1f).

131 Nevertheless, a third group ('CuedRew') was included as a control for satiety effects. This group received  
132 identical rewards to the Deg group, but rather than delivering uncued rewards during the previously blank  
133 trials, these rewards were delivered following a third odor (Odor C). Unlike animals in the Deg group,  
134 animals in the CuedRew group did not decrease anticipatory licking to Odor A. Furthermore, anticipatory  
135 licking, background licking and licking latency were similar to the Cond group (Fig. 1d, e; Extended Data  
136 Fig. 1).

137  $\Delta P(A)$  is 0.25 in the Cued Reward condition, for identical reasoning as the Deg group. This indicates that  
138 the  $\Delta P$  definition of contingency cannot be the sole determinant of conditioned responding (Fig. 1c). This  
139 phenomenon has been previously noted in the behavioral responses in conditioning tasks during  
140 contingency degradation<sup>14,35</sup>. It is not resolved by considering a retrospective definition of contingency.  
141 Consider  $\Delta P_{retro}(A) = P(A + |R) - P(A - |R)$  in both the CuedRew and Deg groups, this quantity is  
142 identical, with Odor A preceding the reward 50% of the time in both conditions.

143 In the subsequent stage of our investigation (Phase 3; 'Recovery 1'), we reinstated the original  
144 conditioning parameters for the Deg group, which increased the contingency back to 0.75 for Odor A,  
145 yielding an immediate recovery of the level of anticipatory licking (Extended Data Fig. 1g).

146 To compare the behavior and neural correlates of contingency, we also introduced an Extinction phase  
147 (Phase 4) to the Deg group. In this phase, no reward was ever delivered following either odor cue. Over  
148 three sessions, anticipatory licking to Odor A gradually waned. Finally, during a second recovery phase

149 (Phase 5; Recovery 2), the anticipatory response to Odor A was effectively reinstated (Extended Data Fig.  
150 1g).

151 Notably, apart from the Extinction phase, the probability of a reward following Odor A was constant at  
152  $P(R|A) = 0.75$  throughout the experiment while behavior changes considerably. Clearly, the contrast  
153 against the probability of reward in the absence of a cue is an important consideration for anticipatory  
154 behaviors, with marked changes during contingency degradation. However, the Cued Reward control  
155 showed it is not as straightforward as the contrast between the absence and presence of a cue.

156

## 157 **Contingency degradation attenuates dopaminergic cue responses**

158 Given the well-documented role of dopamine in associative learning, we sought to characterize the  
159 activity of dopamine neurons in our Pavlovian contingency manipulation task. We monitored axonal  
160 calcium signals of dopamine neurons using a multi-fiber fluorometry system<sup>36</sup> with optical fibers  
161 targeting 6 locations within the ventral striatum (VS), including the nucleus accumbens (NAc, medial and  
162 lateral) and the olfactory tubercle (OT, 4 locations; Fig. 2a, b). Recordings were made only in the Deg and  
163 CuedRew groups, with the final session of Phase 1 used as the within-animal conditioning control. To  
164 ensure similar levels of calcium sensor expression across the six recording locations, we employed a  
165 transgenic approach by crossing a transgenic mouse line expressing the Cre recombinase in dopamine  
166 neurons (DAT-Cre)<sup>37</sup> and a reporter line that expresses a calcium sensor GCaMP6f in a Cre-dependent  
167 manner (Ai148)<sup>38</sup>. Fiber locations were verified using post-mortem histology (see Methods for exclusion  
168 criteria, Fig. 2b). All results presented in the main text are from the lateral nucleus accumbens (INAc),  
169 where TD error-like dopamine signals have been observed most consistently<sup>39</sup>, though the main findings  
170 are consistent across all locations (minimum cosine similarity between any other area and INAc's DA  
171 signals during odor A rewarded trials: 0.92, Extended Data Fig. 2).

172 During Phase 1 (initial conditioning) dopamine axons in INAc initially responded strongly to water and  
173 weakly to Odor A (Fig. 2c, d). As learning progressed, the response to water gradually decreased, while  
174 the response to Odor A increased over the course of 5 sessions ( $t_{13} = 4.81$ ,  $P = 0.0004$ , paired  $t$ -test, cue  
175 response first vs. last session of Phase 1), broadly consistent with previous reports of odor-conditioning  
176 on stochastic rewards<sup>29,40</sup>.

177 During contingency degradation (Deg group, Phase 2), the response to Odor A decreased across sessions  
178 ( $t_8 = -11.50$ ,  $P = 8.4 \times 10^{-6}$ , paired  $t$ -test, cue response, session 6 versus 10) consistent with the changes in  
179 anticipatory licking and other recent reports of dopamine during contingency degradation<sup>10,31,32</sup> (Fig. 2e,  
180 f). However, in the Cued Reward condition (CuedRew group, Phase 2), the response to Odor A did not  
181 decrease compared to the Phase 1 response ( $t_5 = -1.12$ ,  $P = 0.32$ , paired  $t$ -test, cue response first vs. last  
182 session of Phase 2), aligning with the behavioral results but conflicting with the idea that dopamine  
183 neurons encode contingency, at least so far as defined by  $\Delta P$ .

184 In the additional phases (3-5) in the Deg group, dopamine also mirrored behavior: the response to Odor A  
185 quickly recovered in Recovery 1 (Phase 3), decreased during Extinction (Phase 4) and recovered again  
186 during Recovery 2 (Phase 5; Extended Data Fig. 3a). These results show that dopamine cue responses  
187 track the stimulus-outcome contingency in our Pavlovian contingency degradation and extinction  
188 paradigms although they deviated from the contingency in the CuedRew group. Still, in all groups and  
189 phases, dopamine tracked anticipatory licking.

190

## 191 **TD learning models can explain dopamine responses in contingency degradation**

192 In both behavior and dopamine, the responses are not fully explained by contingency: there were  
193 diminished responses during contingency degradation, but not when the additional rewards are cued.  
194 Given the match between dopamine responses and behavior, rather than consider new definitions of  
195 contingency, we sought to test if temporal difference (TD) models, which so far have been highly  
196 successful in accounting for dopamine activity, are able to explain the discrepancies from the contingency  
197 account.

198 Dopamine neurons are thought to convey TD errors, denoted by  $\delta$  and defined by the equation:  $\delta_t = r_t +$   
199  $\gamma V(s_{t+1}) - V(s_t)$ , with  $r_t$  representing reward at time  $t$ ,  $s_t$  representing the state at time  $t$ ,  $V(s_t)$  is the  
200 value at state  $s_t$ , and  $\gamma$  is the temporal discount factor ( $0 < \gamma < 1$ ). Value  $V(s_t)$  is defined as the  
201 expected sum of all future rewards starting from time  $t$ , with each future reward discounted by the factor  
202  $\gamma$  at each time step. The role of the TD error in learning is to iteratively refine the value estimate (Fig. 3a),  
203 ultimately guiding behavior.

204 The response to Odor A differed most between our three test conditions (Conditioning, Degradation,  
205 Cued Reward) and thus our modeling efforts initially focused on explaining these changes. Noting there  
206 was no reward at the time of Odor A, by the definition of TD error, the response to Odor A is  $\gamma V(s_{t+1}) -$   
207  $V(s_t)$ , the difference between the value in the state immediately after Odor A (ISI) and the value in the  
208 state immediately before Odor A (pre-Odor ITI).

209 Previous studies have indicated that the ability of TD learning models to explain dopamine responses and  
210 conditioned behaviors depends critically on what types of state representations the models use<sup>41–44</sup>. We  
211 therefore tested TD learning models (Fig. 3a) equipped with four different types of state representation  
212 (Fig. 3c).

213 In the original application of TD models to dopamine activity, only the interval between a CS and a US,  
214 i.e. inter-stimulus interval (ISI), was considered, and was represented using a ‘complete serial compound’  
215 (CSC) representation, sometimes known as a tapped delay line<sup>25,45</sup>. In this construction, a presentation of  
216 a stimulus triggers a sequential activation of sub-states, each of which represents a time step after the  
217 stimulus (Fig. 3c). At any given time after the stimulus, only one sub-state is active. The value estimate  
218  $\hat{V}(s_t)$  is then computed as the weighted sum of these substates which in CSC reduces to be the weight of  
219 the active substate.

220 While this ISI-only CSC state representation is successful in explaining many properties of the dopamine  
221 response to conditioned stimuli, it fails to predict the result of our experiments. As the ISI period is  
222 identical between conditions and there is no representation of the ITI period, the TD error for Odor A is  
223 unchanged between conditions (Fig. 3f).

224 An extension of this ISI-only model (CSC with ITI states model) models both the ISI and ITI using CSC,  
225 resetting with each odor. While this model predicts a decrease for Degradation, it also predicts a decrease  
226 in the Cued Reward condition (Fig. 3f), conflicting with our results.

227 Rather than representing the ITI with many consecutive states, it is possible to represent it as a single  
228 state. This model, which we term the Cue-Context model, is functionally similar to the previously  
229 developed cue competition model<sup>16–19</sup>. Our Cue-Context model extends the original CSC model with a

230 state that is constantly on (the ‘context’) during both the ISI and ITI (Fig. 3b). This model successfully  
231 predicts the pattern of experimental results we observed, with a decrease in the Odor A response during  
232 Degradation and a smaller decrease during Cued Reward (Fig. 3f). This can be understood as the context  
233 acquiring value, lesser in the Cued Reward condition because the increased value (more rewards) is  
234 attributed to both the context and the new odor. On the other hand, in the Degradation condition, the  
235 increased value is attributed fully to the context. By increasing the context and thus value during the pre-  
236 Odor ITI period, the TD error at Odor A is diminished. While this produces a qualitatively correct pattern  
237 of results, it requires a temporal discount factor that is well below previously reported values<sup>46–49</sup> to  
238 produce the quantitatively correct pattern (Extended Data Fig. 4).

239 We therefore considered whether further information about the experimental design could be used to  
240 refine the state representation. Inspired by previous work showing that dopamine neurons are sensitive to  
241 hidden state inference in a task with stochastically timed rewards<sup>50,51</sup>, we considered a ‘belief state’  
242 representation, a vector of probabilities for each possible hidden state (Belief-State TD model; Fig. 3b):  
243 the ‘Wait’ state, which reflects early ITI (a minimum fixed amount of waiting period in which there is no  
244 chance of an uncued reward or odor being delivered), and the ‘pre-transition’ (Pre) state, in which there is  
245 an imminent chance of reward or odor being delivered. The transition and observation matrix, which are  
246 used to compute the probability of each state, were derived from the experimental settings, assuming a  
247 fixed probability of transition from the Wait to Pre state, modeling a growing anticipation of the next trial  
248 beginning. Using this state representation improved the quantitative accuracy of the model for a given  $\gamma$   
249 versus the Cue-Context, and accurately predicted the experimental data at a value of consistent with  
250 previously reported results<sup>46–49</sup> (Fig 3f., Extended Data Fig. 4).

251 To test which of these two models, Cue-Context or Belief-State, best describes the state representation  
252 driving dopamine responses and behavior, we focused our analysis on the ITI period: whereas the Cue-  
253 Context representation models the ITI as a single, homogenous state (the context), the Belief-State model  
254 captures temporal heterogeneity by modeling it as the gradual transition between two states – capturing a  
255 growing anticipation of the next trial or reward (Fig. 4a).

256 In Pavlovian settings, anticipatory licking (as opposed to consummatory licking) has been used as a  
257 measure of current value – for example, animals will lick more to cues that predict greater rewards<sup>26</sup>.  
258 Odor B provides an opportunity to examine whether the ITI is a heterogenous interval. This is because  
259 Odor B predicts no reward within the current trial (and thus no consummatory licking) but also provides  
260 further information that no odor or uncued reward will be delivered for the length of one trial. In this way,  
261 Odor B provides 100% certainty to the animal that while they are in the ‘Wait’ state. Consistent with this  
262 understanding of the task structure, the delivery of Odor B during the Degradation condition prompted  
263 animals to stop licking. Both the Cue-Context and Belief-State models capture this effect. The crucial  
264 difference is how the lick rate recovers. In the Cue-Context model, ITI value is related to a single state,  
265 which without reward decreases at the rate of  $\alpha$  (learning rate). In the Belief-State model, value  
266 continually increases (Fig. 4c) across the entire ITI, as the increased belief that the next trial is imminent  
267 increases continuously. We find that the lick pattern following Odor B matches the pattern of the Belief-  
268 State model and not the Cue-Context model, with a sudden decrease in licking followed by a gradual  
269 increase in the Degradation condition that is unrelated to the ISI length (Fig. 4c, summarized in Fig. 4d).

270 This pattern of licking behavior also suggests that the animals do not develop more complex models of  
271 timing. Odor B predicts approximately ten seconds with no reward. An ideal agent would not lick during

272 this time, waiting until the transition to an uncued reward is possible. The mice instead resume licking  
273 within 2-3 seconds of Odor B delivery, with the lick rate increasing over several seconds.

274 While the value from the Belief-State model explains the time course of licking following Odor B, this  
275 account does not, by itself, explain the decrease in anticipatory licking in response to Odor A (Fig. 1d).  
276 This decreased responding is a consistent feature of contingency degradation<sup>10,34</sup>. We show, in Extended  
277 Data Fig 5, that if licks carry a small effort cost and licks are distributed according to relative value, then  
278 the Belief-State model can account for the increased licking in the pre-odor period and decreased licking  
279 during the ISI.

280 Having shown that the lick rate is explained by the changing value in the Belief-State model, we wished  
281 to test whether this could be used to explain trial-by-trial variance in the dopamine response. Continuing  
282 with the assumption that licking is a moment-by-moment measure of value, our Belief-State model  
283 predicts there should be an inverse correlation between the pre-odor lick rate and the Odor A dopamine  
284 response. To test this, we correlated the number of licks in the two seconds before the cue to the Odor A  
285 response on a trial-by-trial basis, regressing a linear model independently for each mouse (pooling the last  
286 two sessions of each condition under the assumption that the task was well-learned in these sessions).

287 Only in the Degradation condition was there a significant negative correlation between the pre-odor lick-  
288 rate and the Odor A dopamine response for the population (Fig. 4f, g). This can be explained by the  
289 Belief-State model, as ITI value varies depending on the length of the ITI – with each timestep, there is an  
290 increasing belief they are in the ‘Pre’ state, with the current value estimate updating to reflect that. The  
291 data cannot be explained by the Cue-Context model, in which ITI value is fixed (Fig. 4g). The modeling  
292 suggests that the lack of a significant trend in the remaining two conditions is due to the lower variance in  
293 value in the pre-cue period, with an average of  $0.28 \pm 0.87$  and  $0.46 \pm 1.11$  licks (mean  $\pm$  s.d.) in the 2  
294 second pre-cue period in Conditioning and Cued Reward respectively (versus  $1.51 \pm 1.52$  in  
295 Degradation), leading to underpowered analysis.

296 In summary, the ITI state representation is essential to explaining the relative effects of contingency  
297 degradation and additional cued rewards on the Odor A response. Complex ITI representations, such as  
298 CSC, are inefficient, whereas modeling it as a single state (Cue-Context), does not capture the  
299 heterogeneity of the ITI. Our Belief-State model, representing the ITI using two states, is sufficient to  
300 explain the experimental results.

301

## 302 **Additional aspects of dopamine responses and model predictions**

303 Having identified a sufficient model for explaining our contingency degradation results, we next  
304 examined how well this model matched additional experimental results. Figure 5a visualizes the value as  
305 predicted by the Belief-State model across four conditions tested (Conditioning, including Recovery;  
306 Degradation, Cued Reward and Extinction). In the Odor A rewarded trial, the value during the ISI  
307 remained unchanged in the first three conditions, and significantly decreased in Extinction, closely  
308 mirroring the (prospective) reward probability  $P(R|A)$ . For the reasons discussed above, the pre-ISI  
309 period, reflecting the pre-transition state (‘Pre’), showed a modest increase in the Cued Reward case and a  
310 significant rise in the Degradation condition. The TD errors upon Odor A presentation, reflective of the  
311 difference in value between Pre state and the first ISI substate, diminished in both Degradation and  
312 Extinction. In both these conditions, contingency is reduced by increasing  $P(R|A-)$  and decreasing



313  $P(R|A)$ , respectively<sup>52</sup>. Notably, our model suggested two distinct mechanisms underlying these two  
314 processes: an increase in Pre state value in Degradation and a decrease in ISI value in Extinction (Fig. 5c).  
315 Our Belief-State TD learning model matched the experimental results well (Fig. 5b, d), including the  
316 Extinction data.

317 The model predicts another distinct difference between degradation and extinction: degradation affects  
318 TD error for all cues due to changes in the shared Pre state value, while extinction impacts only the  
319 specific cue undergoing extinction. Accordingly, we examined the Odor B trials. In the Belief-State  
320 model, Odor B is a transition from the Pre to Wait state, and thus the TD error is the difference between  
321 these two state values. We expected the most negative response in the Deg group, owing to a higher Pre  
322 state value, and relatively unchanged ‘Wait’ value. We also expected an unchanged response in  
323 Extinction in comparison to Conditioning. Experimentally, the response to Odor B was biphasic,  
324 featuring an initial positive response followed by a later negative response. Such a biphasic response has  
325 been previously noted, with general agreement that the second phase is correlated with value<sup>53</sup>. By  
326 quantifying the later response (250ms-1s window), there was a close match between the model prediction  
327 and the data for Odor B responses (Fig. 5e, f).

328 The Belief-State model shows that TD errors at reward omission are based on the difference between the  
329 final ISI substate and Wait state values. The Wait state value, generally lower than the Pre state value, has  
330 minor changes across conditions. This results in consistent TD errors at reward omission across  
331 Conditioning, Degradation, and Cued Reward conditions due to similar ISI values, but a significant  
332 reduction in Extinction due to a lower ISI value, which closely aligned with the experimental results  
333 (Extended Data Fig. 6). TD errors at predicted rewards, reflecting the difference between actual reward  
334 and ISI values, exhibit minimal changes across Conditioning, Degradation and Cued Reward conditions,  
335 which is also consistent with the data.

336 In total, the above results indicate that the TD model with proper task states can effectively recapitulate  
337 nearly all aspects of phasic dopamine responses across various trial types and task events.

338

### 339 **Recurrent neural networks that learn to predict values through TD learning can explain** 340 **dopamine responses**

341 The models discussed above, while effective, are ‘hand-crafted’ and tuned to our particular task setting.  
342 While there is evidence that dopamine neurons rely on belief-state inference in computing TD error<sup>50,51,54</sup>,  
343 the question of how the brain learns such a state-space is less well understood. Previous work has shown  
344 that RNNs, trained to estimate value directly from observations (‘value-RNNs’), develop belief-like  
345 representations despite not being explicitly trained to do so<sup>55</sup>. This approach substitutes hand-crafted  
346 states for an RNN that is simply given the same odor and reward observations as the animal (Fig. 6a).

347 Here, we applied the same value-RNN to our contingency manipulation experiments. We generated  
348 training sets, consisting only of the odor and reward timings, that matched the three conditions. The  
349 RNNs were first trained on the Phase 1 Conditioning task and then either on the Phase 2 contingency  
350 Degradation or Cued Reward conditions (Fig. 6b). Several RNNs were trained with different numbers of  
351 hidden units, from 5 to 50.

352 The trained RNNs closely matched the experimental results (example 50 unit RNN presented in Fig. 6c).  
353 Like the TD models used in the above section, the decrease in Odor A response is explained by an  
354 increase in the value during the ITI period, not a shift in the value during the ISI (example Fig. 6d).

355 We were interested in understanding the inferred state spaces used by the RNN models. To visualize this,  
356 we applied canonical correlation analysis (CCA)<sup>56,57</sup> to align the activity of the hidden units between the  
357 RNNs for each condition for all conditions.

358 In all conditions, without any stimuli, the RNN's activity will decay to a fixed point (here plotted as the  
359 origin, Extended Data – Video 1). This can be understood as the Pre-transition state. In all conditions, the  
360 Odor A trajectory is similar, indicating a shared representation of the ISI period (Fig. 6e). Furthermore, in  
361 the Cued Reward condition, the Odor C trajectory is nearly identical to that of Odor A, potentially  
362 reflecting generalization. In the Degradation condition, delivering Odor B causes a trajectory that is  
363 significantly longer than the other two conditions, potentially corresponding to the Wait state.

364 To compare the state space of the value-RNN to the Belief-State model, we calculated the beliefs at each  
365 given time point in the simulated experiment and used a linear regression to relate the hidden unit activity.  
366 As previously noted<sup>55</sup>, the unit activity became more belief-like with more hidden units (Fig. 6f). Notably,  
367 the regression performance, as quantified by  $R^2$  (see Methods), was higher for the Degradation condition  
368 at each hidden layer size. This is explained by better performance on the Wait state (Fig. 6f, right panel).  
369 As evident in the visualized activity in the state spaces, the RNNs trained on the Degradation condition  
370 developed distinct trajectories in the ITI compared to the other two conditions (Fig. 6g), taking a longer  
371 period of time to return to the fixed ITI point and following a similar pattern regardless of the particular  
372 trial type. In all RNNs that successfully predicted degradation effect, the Wait state readout had a  
373 minimum performance of  $R^2 = 0.57$ . This suggests that it is the delivery of rewards during the ITI that  
374 reshapes the state space to be heterogeneous, while in the other conditions this is not necessary and thus  
375 the ITI has a relatively fixed state space representation. That the RNN can learn a belief-like  
376 representation from limited information, using only the TD error as feedback, suggests a generalized  
377 method by which the brain can construct state spaces using TD algorithms.

378

### 379 **A retrospective learning model, ANCCR, cannot explain the dopamine responses**

380 While our analysis using the TD learning models with explicit state representations and the value-RNNs  
381 suggest that TD learning models are sufficient to explain our experimental results, we have not yet  
382 considered whether alternative definitions of contingency would also provide an account of our results.  
383 The ANCCR (adjusted net contingency for causal relations) is a recently described new model, proposed  
384 as an alternative account of the TD explanation of dopamine activity (Fig. 7a)<sup>31</sup>. The authors have  
385 previously shown that this model can account for contingency degradation<sup>31,32</sup> and suggested that TD  
386 accounts could not.

387 ANCCR builds upon the authors' previous observation that the retrospective information ('which cues  
388 precede reward?') can be used to explain animal behavior previously unexplained by prospective  
389 accounts<sup>33</sup>. Accordingly, the ANCCR model begins with the calculation of the retrospective contingency,  
390 using eligibility traces as a principled method to compute contingency in continuous time, rather from  
391 trial-by-trial probabilities. At the time of a reward (or a 'meaningful event'), the difference between the  
392 eligibility of cues at the time of reward and the average cue eligibility is computed. This generalizes the

393 trial-based definition of  $\Delta P_{retro}(A)$  to continuous time. From this retrospective contingency and the  
394 average event rates, the model proposes the prospective contingencies are inferred using a Bayes-like  
395 computation.

396 From the prospective and retrospective contingencies, a weighted sum ('net') contingency is calculated  
397 for all pairs of events. This map can be used to calculate the change in expectation of reward for a given  
398 event, considering other explanations. It is this 'adjusted net contingency' that ANCCR proposes is  
399 represented in the dopamine signal.

400 To test the ANCCR model, we used the authors' published code to model the same 25 simulated  
401 experiments used in our TD modeling. For this experiment, ANCCR has 12 parameters, at least 6 of  
402 which have a significant impact on the modeled response. We first tried using the parameters published in  
403 Garr et al. (2023) and Jeong et al. (2022); we present results using the Garr parameters because they are  
404 closer to our experimental results. While the ANCCR model accurately predicted a decreased response  
405 for Odor A during contingency degradation, it predicted a similar response in the Cued Reward condition,  
406 conflicting with the experimental results (Fig. 7b). We varied the parameter ( $w$ ) which controls the  
407 relative amount of retrospective and prospective information used to calculate the contingency. This  
408 parameter controlled the relative size of the decrease, while sensitive to parameter choice governing the  
409 eligibility decay rate and learning rates; in general the halving of  $P(A|R)$  produces a large decrease in the  
410 retrospective contingency,  $P(A|R) - P(A)$ , whereas the increase in  $P(R)$  slightly decreases the  
411 prospective contingency,  $P(R|A) - P(R)$ .

412 We next considered whether this was a problem of parameter selection, and therefore simulated the first 5  
413 virtual experiments for the parameter search space used in Garr et al. (2023), trying a total of 21,000  
414 parameter combinations, including those in the two previous studies<sup>31,32</sup> (indicated as 1, 2 and 3). Fig. 7c  
415 plots the Odor A dopamine response in the Degradation and Cued Reward case for each of these  
416 combinations, normalized by the response during Conditioning. No parameter combination predicted the  
417 correct pattern of experimental results, quantitatively or qualitatively (Fig. 7c).

418

419

## 420 Discussion

421 Here we examined behaviors and VS dopamine signals in a Pavlovian contingency degradation paradigm,  
422 including a pivotal control. Our results show that dopamine cue responses, like behavioral conditioned  
423 responses, were attenuated when stimulus-outcome contingency was degraded by the uncued delivery of  
424 additional rewards. Crucially, neither dopamine signals nor conditioned responses were affected in a  
425 control condition in which the delivery of additional rewards were cued by a different stimulus despite a  
426 similar number of rewards being administered. Our findings not only demonstrate that the above results  
427 were not due to satiety, but also provide key insights into possible mechanisms underlying contingency  
428 degradation.

429 Contrary to claims from previous studies<sup>31,32</sup>, our modeling showed that many aspects of dopamine  
430 responses can be comprehensively explained by TD learning models, if the model is equipped with proper  
431 state representations reflecting the task structure. These TD learning models also readily explained

432 dopamine cue responses in the control condition with cued rewards – results which strongly violated the  
433 predictions of a contingency-based retrospective causal learning model (ANCCR) and the  $\Delta P$  definition  
434 of contingency. The results indicate that dopamine signals, as well as conditioned responses, primarily  
435 reflected the prospective, but not retrospective, stimulus-outcome relations. Rather than discarding the  
436 notion of contingency altogether, we propose that these results point toward a novel definition of  
437 contingency grounded in the TD learning framework. These results bear significant implications for the  
438 theory of associative learning and the nature of dopamine signals, which help resolve some of previously  
439 unresolved controversies.

#### 440 **TD learning model as a model of associative learning**

441 Historically, Pavlovian contingency degradation paradigms have played a pivotal role in the development  
442 of animal learning theories<sup>3,16</sup>, yet the exact underlying mechanisms remain to be determined<sup>17,19</sup>. Here we  
443 show that the effect of contingency manipulations, both on behavior and dopamine responses, can be  
444 explained by TD learning models. As our systematic investigation revealed, the failure of previous efforts  
445 to explaining contingency degradation with TD learning models is due to the use of inappropriate state  
446 representations, either not considering the ITI period at all, or modeling it in a simplistic way. We show  
447 two types of TD learning models that explain the basic behavioral and dopamine results. The first model  
448 (Cue-Context model) uses a contextual stimulus as one of the states that continuously exists throughout  
449 the task period, which is equivalent to the cue-competition model traditionally considered in the animal  
450 learning theory literature<sup>16,17,19</sup>. The second model (Belief-State model) explicitly models the task  
451 structure as transitions across the ISI and the two ITI states (“wait” and “pre-transition”), and TD learning  
452 operates on beliefs (posterior probabilities) over these discrete states.

453 In both models, the reduction in dopamine cue responses occurs due to an increase in the value preceding  
454 a cue presentation, which decreases the *change* in value (reward expectation) induced by the cue, rather  
455 than due to a decrease in the absolute level of the value induced by the cue. This raises the question of  
456 why cue-induced anticipatory licking is reduced during contingency degradation. We provide a potential  
457 mechanism: the animal distributes anticipatory behavior depending on the relative values across different  
458 states.

459 Our results favor the Belief-State model over the Cue-Context model; both the dopamine and behavioral  
460 data were better explained by the Belief-State model. One could argue that it is unclear whether the  
461 animal can learn “sophisticated” state representations such as those used in our Belief-State model. In  
462 support of a Belief-State TD learning model, our analysis of anticipatory licking indicated that the reward  
463 expectation was modulated in a manner intricately linked to different task states: the ISI, wait, and pre-  
464 transition states. Furthermore, we show that recurrent neural networks, trained to predict values (value-  
465 RNNs), acquired the activity patterns that can be seen as representing beliefs, merely from observations,  
466 without explicitly instructed to develop such representations, similar to our previous work using different  
467 behavioral tasks<sup>55</sup>. Critically, when trained on contingency degradation sessions, the value-RNNs  
468 developed more heterogenous representations of the ITI, capturing the same phenomenon as the Belief-  
469 State model.

470 It has been shown that TD learning models can explain a wealth of phenomena studied in the animal  
471 learning theory literature<sup>30,58</sup>. The present study adds to this list Pavlovian contingency degradation – a

472 classic phenomenon long studied in psychology and now in neurobiology. These results indicate that TD  
473 learning models provide a foundation with which to understand associative learning while the RNN-based  
474 approach provides a principled way to apply TD learning with minimal assumptions about state  
475 representations.

#### 476 **State representations as population activity dynamics**

477 In RL, the “state” is a critical component which represents the set of observable and inferred variables  
478 necessary to compute value and policy. The artifice of the state representations used in neurobiological  
479 RL modeling has been criticized<sup>59</sup>. For instance, it is implausible to have separate sets of neurons  
480 activated sequentially (i.e. CSCs) for separate cues, particularly if they are to completely tile the ITI, as in  
481 the CSC with ITI states model<sup>59</sup>. Furthermore, states are often defined within each “trial”; how can states  
482 be defined when there are no obvious trial structures<sup>59</sup>? The success of value-RNNs in replicating aspects  
483 of dopamine signals and the acquisition of belief-like representations provides two crucial insights into  
484 how biological circuits may represent states.

485 First, the recent successes of RL on complex machine learning tasks, containing many stimuli and often  
486 lacking obvious trial structure, suggests that it is possible to achieve high performance with standard RL  
487 techniques<sup>22</sup>. A key ingredient lies in the use of neural networks capable of autonomously learning  
488 representations appropriate for specific tasks. Our results with value-RNNs agree with this observation.  
489 As shown in our previous work<sup>55</sup> and in the present work, value-RNNs have a stable fixed point  
490 (attractor) corresponding to the ITI state (the Pre-transition state in our Belief-State model). The ITI state  
491 is thus an emergent property of training to predict reward. Furthermore, different stimuli induce stimulus-  
492 specific trajectories in the population activity state space. We found a close correspondence between  
493 population dynamics and the hand-crafted states assumed in Belief-State TD learning models. These  
494 results indicate that the population activity patterns in a network, including attractors and stimulus-  
495 specific trajectories, represent distinct states such as those assumed in our Belief-State TD learning  
496 models. Although the activity representing different trajectories likely involves overlapping sets of  
497 neurons, they can be trained to compute values properly by adjusting downstream synaptic weights (as  
498 long as the activity patterns for different states are discriminable).

499 Second, while TD learning models with hand-crafted state representations help develop conceptual  
500 understanding, the RNN-based approach can provide insights into how hand-crafted state representations  
501 could be implemented in neural networks. In the future, it is of great interest to examine whether neural  
502 activity in the brain exhibits patterns of activity predicted by the value-RNN models. The prefrontal  
503 cortex is a strong candidate area, receiving dopaminergic innervation from the VTA necessary for  
504 appropriate adaptation to contingency degradation in instrumental conditioning<sup>60</sup>. However other areas,  
505 such as the hippocampus, also contribute task-relevant information during degradation to the prefrontal  
506 cortex<sup>61</sup> and neural network modeling approaches that reflect the brain’s functional organization (e.g. <sup>62</sup>)  
507 might provide more insight than our model which treats the state-machinery of the brain as a single  
508 recurrent neural network.

#### 509 **Limitations of the ANCCR model as a model of associative learning and dopamine**

510 The present study unveiled limitations of the recently proposed causal learning model, ANCCR<sup>31,32</sup>. The  
511 Degradation and Cued Reward conditions are minimally different and thus provide a strong test of the  
512 algorithm design. Our results indicate that the ANCCR model fails to explain the observed results despite  
513 our extensive examination of its parameter space. Crucially, the ANCCR model suffers from the same  
514 flaw as the  $\Delta P$  definition of contingency. While extending the definition to continuous time and  
515 considering multiple cues, ANCCR still calculates contingency by subtracting the background event rate,  
516 losing information, precluding it from attributing increased value to the background in the same manner  
517 as the TD models. Given the similarity in event rate between the conditions, the retrospective  
518 representations (and average eligibility trace) remain similar, with the computed retrospective  
519 contingencies in the Degradation condition being a subset of the Cued Reward contingencies (Fig 7d).  
520 This explains why ANCCR predictions are similar for the two conditions independent of the parameter  
521 choice, as the rest of the model depends on this computed retrospective contingency as input. Thus, the  
522 failure of the ANCCR to explain the Cued Reward condition reflects the fundamental construction of the  
523 ANCCR model.

524 The failure of the ANCCR model here does not exclude some of the interesting ideas integrated into  
525 ANCCR, including how it uses retrospective information to learn the state space. Rather, it is its reliance  
526 on contingency that constitutes its core deficit. Other theoretical work has considered how TD algorithms  
527 that consider retrospective information may enhance learning performance without explicitly invoking  
528 contingency. A recent report<sup>32</sup> demonstrated that ANCCR is able to explain the dopamine response in  
529 outcome-selective contingency degradation. This is a result of the multidimensional tracking of cue-  
530 outcome contingencies in ANCCR. We show that both the Belief-State model and the value-RNN, if  
531 trained on each reward separately and with total value taken to be the absolute difference of the two  
532 separate values, successfully predicts the experimental results of Garr et al. (2023) (Extended Data Fig.  
533 7). A similar approach using “multi-threaded predictive models” was used to successfully explain  
534 dopamine data in a different multi-outcome task<sup>44</sup>. While this proposal leaves open questions about how  
535 such abstract state representation is implemented biologically (the same being true for ANCCR), it does  
536 demonstrate that more complex contingency manipulations can still be explained by TD models. In fact,  
537 recent studies have provided evidence for heterogeneous responses to different types of rewards in  
538 dopamine neurons<sup>63–65</sup>. While further evidence is required to solidify this understanding, the provisional  
539 assumption of multiple value channels shows how TD models for multiple outcomes can potentially be  
540 achieved in neural circuitry by concurrently running parallel circuits.

#### 541 **TD error, contingency and causal inference**

542 Learning predictive or causal relationships requires properly assigning credits to those events that are  
543 responsible for the outcomes observed. A key to this process is considering counterfactuals<sup>66</sup> – would a  
544 particular outcome occur had I not seen that cue, or had I not taken that action? In the present study, we  
545 show that TD learning models with ITI representations learn and predict value in the time before cue  
546 presentation. The cue-associated TD error is then calculated as the difference in value in the presence and  
547 absence of that cue. Consequently, computation of TD errors effectively subtracts the prediction of value  
548 in the absence of the cue – i.e. the counterfactual prediction. More generally, the computation of TD error  
549 or its variants can be seen as subtracting out counterfactuals. In a class of RL algorithms commonly used  
550 in artificial intelligence applications (advantage actor-critic algorithms), the actor decides which action to  
551 take for a given state and the critic evaluates the action by computing the advantage function, defined as:

552 
$$A(s_t, a_t) = Q(s_t, a_t) - V(s_t)$$

553 where  $Q(s_t, a_t)$  is the state-action value function<sup>67,68</sup>. If this is taken to be the immediate reward of the  
554 action plus the expected return of the new state,  $Q(s_t, a_t) = r_t + \gamma V(s_{t+1})$  then the advantage function  
555 can be approximated by the TD error  $A(s_t, a_t) = Q(s_t, a_t) - V(s_t) = (r_t + \gamma V(s_{t+1})) - V(s_t) = E[\delta_t]$   
556 (ref. <sup>69,70</sup>).

557 As discussed in recent work<sup>70</sup>, in fully observable environments without confounds, the advantage  
558 function is exactly equivalent to the Neyman-Rubin definition of causal effect of an action: the difference  
559 in outcomes given an action versus outcomes otherwise. In this context, the definitions of causality,  
560 contingency and TD error align – all emphasizing the consideration of counterfactual prediction: that is,  
561 the difference between potential future outcomes (following action) and the alternative when the action is  
562 not taken. TD error can therefore be both a measure of contingency and useful in establishing causal  
563 relationships, without invoking retrospective computations.

564 TD errors improve over the ANCCR and  $\Delta P$  definitions because the comparison to the reward probability  
565 of US given CS is not simply the reward probability given absence of CS, but to  $V(s)$ , which is the  $\gamma$ -  
566 discounted sum of all future rewards given the current state, with the state encapsulating all  
567 environmental information. As demonstrated by our modeling, the heterogeneous state representation  
568 during the absence of events (the ITI) is critical to the accuracy of our models to match the experimental  
569 data.

570 While these relationships between TD error and contingency hold in fully observable environment, our  
571 value-RNN approach may extend these results to more complex/realistic environments. Veitch et al.  
572 (2019) has demonstrated network embeddings, like our value-RNN, can reduce the problem of inferring  
573 causality to a problem of predicting outcomes<sup>71</sup>. These networks do not require full knowledge of the  
574 environment to succeed but rather learn to extract sufficient information to establish causality. Ultimately,  
575 TD error could provide pivotal signals for contingency – the essential quantity for causal inference.

## 576 **Conclusions**

577 Our results indicate that TD learning models can explain contingency degradation – a phenomenon that  
578 was thought to be difficult to explain based on TD learning<sup>31,32,72</sup>. The Belief-State TD model that we  
579 used here is “model-free” in the sense that the values are “cached” to each state based on direct  
580 experiences, although these states reflect the animal’s knowledge of the transition structure between states  
581 which can be regarded as a “world model”<sup>41,43,73</sup>. This suggests that the distinction between “model-free”  
582 and “model-based” mechanisms is not as hard-lined as often assumed. The sensitivity to contingency  
583 degradation in instrumental behaviors has been used to support the behavior being goal-directed or  
584 model-based. Yet, the same type of Belief-State TD model can, in principle, be applied to explain such an  
585 effect. In any case, further biological investigations will be needed to constrain mechanisms linking  
586 behavior and contingency – the critical variable thought to underlie learning predictive and/or causal  
587 relationships. The experimental results and models presented in this study would aid such efforts.

## 588 **Methods**

589

### 590 **Animals**

591 A total of 31 mice were used. 18 wildtype mice (8 males and 10 females) at 3–6 months of age were used  
592 to collect only behavioral data. For fiber photometry experiments, 13 double transgenic mice resulting  
593 from the crossing of DAT-Cre (Slc6a3tm1.1(cre)Bkmm; Jackson Laboratory, 006660)<sup>37</sup> with Ai148D  
594 (B6.Cg-Igs7tm148.1(tetO-GCaMP6f,CAG-tTA2)Hze/J; Jackson Laboratory, 030328)<sup>38</sup> (DAT::cre x  
595 Ai148, 7 males and 6 females) at 3–6 months of age were used. Mice were housed on a 12 hr /12 hr  
596 dark/light cycle. Ambient temperature was kept at  $75 \pm 5$  °F and humidity below 50%. All procedures  
597 were performed in accordance with the National Institutes of Health Guide for the Care and Use of  
598 Laboratory Animals and approved by the Harvard Animal Care and Use Committee.

599

### 600 **Surgery**

601 Mice used for fiber photometry recordings underwent a single surgery to implant a multifiber cannula and  
602 a head fixation plate 2-3 weeks prior to the beginning of the behavioral experiment. All surgeries were  
603 performed under aseptic conditions. Briefly, mice were anesthetized with an intraperitoneal injection of a  
604 mixture of xylazine (10 mg/kg) and ketamine (80 mg/kg) and placed in a stereotaxic apparatus in a flat  
605 skull position. During surgery, the bone above the Ventral Striatum area was removed using a high-speed  
606 drill. A custom multifiber cannula (6 fibers, 200  $\mu$ m core diameter, 0.37 NA, Doric Lenses) was lowered  
607 over the course of 10 min to target 6 subregions in the Ventral Striatum. The regions' coordinates relative  
608 to Bregma (in mm) were: Lateral nucleus accumbens (lNAc, AP:1.42, ML:1.5, DV:-4.5); Medial NAc  
609 (mNAc, AP:1.42, ML:1, DV:-4.5); anterior lateral olfactory tubercle (alOT, AP:1.62, ML:1.3, DV:-4.8);  
610 posterior lateral OT (plOT, AP:1.00, ML:1.3, DV:-5.0); anterior medial OT (amOT, AP:1.62, ML:0.8,  
611 DV:-4.8); posterior medial OT (pmOT, AP:1.00, ML:0.8, DV:-5.0). Dental cement (MetaBond, Parkell)  
612 was then used to secure the implant and custom headplate and to cover the skull. Mice were singly housed  
613 after surgery and post-operative analgesia was administered for 3 days (buprenorphine ER-LAB 0.5  
614 mg/ml). Mice used for behavioral training underwent a similar surgical process, but only a head fixation  
615 plate was implanted.

616

### 617 **Behavioral training**

618 After recovery from headplate-implantation surgery, animals were given ad libitum access to food and  
619 water for 1 week. Before experiments and throughout the duration of the experiments, mice were water  
620 restricted to reach 85–90% of their initial body weight and provided approximately 1–1.5 mL water per  
621 day in order to maintain the desired weight and were handled every day. Mice were habituated to head  
622 fixation and drinking from a waterspout 2-3 days prior to the first training session. All tasks were run on a  
623 custom-designed head-fixed behavior set-up, with software written in MATLAB and hardware control  
624 achieved using a BPod state machine (1027, Sanworks). A mouse lickometer (1020, Sanworks) was used  
625 to measure licking as infra-red beam breaks. The water valve (LHDA1233115H, The Lee Company) was  
626 calibrated, and a custom-made olfactometer was used for odor delivery. The odor valves  
627 (LHDA1221111H, The Lee Company) were controlled by a valve driver module (1015, Sanworks) and a



628 valve mount manifold (LFMX0510528B, The Lee Company). All components were controlled through  
629 the Bpod state machine. Odors (1-hexanol, d-limonene, and ethyl butyrate, Sigma-Aldrich) were diluted  
630 in mineral oil (Sigma-Aldrich) 1:10, and 30  $\mu$ L of each diluted odor was placed on a syringe filter (2.7-  
631  $\mu$ m pore size, 6823-1327, GE Healthcare). Odorized air was further diluted with filtered air by 1:8 to  
632 produce a 1 liter/min total flow rate. The identity of the rewarded and non-rewarded odors were  
633 randomized for each animal.

634 In Conditioning sessions, there are three types of trials: (1) trials of Odor A (40% of all trials) associated  
635 with a 75% chance of water delivery after a fixed delay (2.5 s), (2) trials of unrewarded Odor B (20% of  
636 all trials) as control to ensure that the animals learned the task, and (3) background trials (40% of all  
637 trials) without odor presentation. Rewarded odor A trials consists of 2s pre-cue period, 1s Odor A  
638 presentation, 2.5s fixed delay prior to a 9  $\mu$ L water reward and 8s post-reward period. Unrewarded Odor  
639 B trials consist of a 2s pre-cue period, 1s Odor B presentation, and 10.5s post-odor period. Background  
640 trials in the Conditioning phase span a 13.5s eventless period. Trial type was drawn pseudo-randomly  
641 from a scrambled array of trial types maintaining a constant trial type proportion. Inter-trial-intervals (ITI)  
642 following the post-reward period were drawn from an exponential distribution (mean: 2s).

643 Learning was assessed principally by anticipatory licking detected at the waterspout for each trial type,  
644 with mice performing 100-160 trials per session until they reach an asymptotic task performance,  
645 typically after 5 sessions.

646 After the Conditioning phase, the mice were divided into three groups to undergo different conditions:  
647 Degradation (Deg group), Cued Reward (CuedRew group), and Conditioning (Cond group). The Deg  
648 group experienced contingency decrease during the Degradation phase. In the Degradation phase, Odor A  
649 still delivers water reward with 75% probability, and Odor B remains unrewarded. The difference was the  
650 introduction of uncued rewards (9  $\mu$ L water) in 75% of background trials to diminish the contingency.  
651 Animals underwent 5 sessions, each with 100-160 trials, to adapt their conditioned and neural responses  
652 to the new contingency. Degradation changed the cue value relative to the background trial but did not  
653 impact the reward identity, reward magnitude, or delay to/probability of expected reward.

654 The CuedRew group was included to account for potential satiety effects due to the extra rewards the Deg  
655 group mice received in the background trials. Unlike the Deg group, the CuedRew group's background  
656 trials were substituted with rewarded Odor C trials, where mice received additional rewards signaled by a  
657 distinct odor (Odor C). Rewarded odor C trials have the same trial structure as the rewarded odor A trials  
658 and animals were given 5 sessions, with 100-160 trials each, to adapt their conditioned response and  
659 neural responses to this manipulation.

660 The Cond group proceeded with an additional five Conditioning sessions, keeping the trial structure and  
661 parameters unchanged as in the Conditioning phase.

662 Post-degradation: eight mice were randomly chosen from the Deg group for the reinstatement phase,  
663 replicating the initial Conditioning conditions. After three reinstatement sessions, once the animals'  
664 performance rebounded to pre-degradation levels, we initiated the extinction process. This involved the  
665 delivery of both odors A and B without rewards, effectively extinguishing the cue-reward pairing. To  
666 mitigate the likelihood of animals generating a new state to account for the sudden reward absence, a  
667 shorter reinstatement session was conducted prior to the Extinction session on the extinction day.  
668 Extinction was conducted over three days, each day featuring 100-160 trials. After Extinction, a second

669 reinstatement session was implemented, re-introducing the 75% reward contingency for odor A. All eight  
670 animals resumed anticipatory licking within ten trials during this reinstatement.

671

## 672 **Fiber photometry**

673 Fiber photometry allows for recording of the activity of genetically defined neural populations in mice by  
674 expressing a genetically encoded calcium indicator and chronically implanting optic fiber(s). The fiber  
675 photometry experiment was performed using a bundle-imaging fiber photometry setup  
676 (BFMC6\_LED(410-420)\_LED(460-490)\_CAM(500-550)\_LED(555-570)\_CAM(580-680)\_FC, Doric  
677 Lenses) that collected the fluorescence from a flexible optic fiber bundle (HDP(19)\_200/245/LWMJ-  
678 0.37\_2.0m\_FCM-HDC(19), Doric Lenses) connected to a custom multifiber cannula containing 6 fibers  
679 with 200- $\mu$ m core diameter implanted during surgery. This system allowed chronic, stable, minimally  
680 disruptive access to deep brain regions by imaging the top of the patch cord fiber bundle that was attached  
681 to the implant. Interleaved delivery 473 nm excitation light and 405 nm isosbestic light (using LEDs from  
682 Doric Lenses) allows for independent collection of calcium-bound and calcium-free GCaMP fluorescence  
683 emission in two CMOS cameras. The effective acquisition rate for GCaMP and isosbestic emissions was  
684 20Hz. The signal was recorded during each session when the animals were performing the task.  
685 Recording sites which had weak or no viral expression or signal were excluded from analysis.

686 The global change of signals within a session was corrected by a linear fitting of dopamine signals  
687 (473nm channel) using signals in the isosbestic channel during ITI and subtracting the fitted line from  
688 dopamine signals in the whole session. The baseline activity for each trial ( $F_{0 \text{ each}}$ ) was calculated by  
689 averaging activity in the pre-stimulus period between -2 to 0 seconds before an odor onset for odor trials  
690 or water onset for uncued reward trials. Z-score was calculated as  $(F - F_{0 \text{ each}})/\text{STD\_ITI}$  with STD\_ITI the  
691 standard deviation of the signal during the ITI.

692 To quantify Odor A responses, we looked for ‘peak responses’ by finding the point with the maximum  
693 absolute value during the 1-s window following the stimulus onset in each trial. To quantify Odor B  
694 responses, we measured area under curve by summing the value during the 250 ms to 1s window  
695 following the stimulus onset in each trial. This is to separate out the initial activation (odor response) that  
696 we consistently observed, and which may carry salience or surprise information independent of value. To  
697 quantify reward responses, we looked for ‘peak responses’ by finding the point with the maximum  
698 absolute value during the 1-s window following the reward onset in each trial. To quantify reward  
699 omission responses, we looked for area under curve by summing the value during the 0-1.5s window  
700 following the reward omission in each trial.

701

## 702 **Histology**

703 To verify the optical fiber placement and GCaMP expression, mice were deeply anesthetized with an  
704 overdose of ketamine-medetomidine, and perfused transcardially with 0.9% saline followed by 4%  
705 paraformaldehyde (PFA) in PBS at the end of all experiments. Brains were removed from the skull and  
706 stored in PFA overnight followed by 0.9% saline for 48 hours. Coronal sections were cut using a  
707 vibratome (Leica VT1000S). Brain sections were imaged using fluorescent microscopy (AxioScan slide  
708 scanner, Zeiss) to confirm GCaMP expression and the location of fiber tips. Brain section images were

709 matched and overlaid with the Paxinos and Franklin Mouse Brain Atlas cross-sections to identify imaging  
710 location.

711

## 712 **Computational Modeling**

### 713 Simulated Experiments

714 To compare the various models, we generated 25 simulated experiments of Cond, Deg and CuedRew  
715 groups, matching trial statistics to the experimental settings, but increasing the number of trials to 4,000  
716 in each phase to allow to test for steady-state response in both these TD simulations and the ANCCR  
717 simulations. We then calculated the state representation of the simulated experiments for each of four  
718 state representations (CSC with and without ITI states, Context-TD, Belief-State model, detailed below)  
719 and ran the TD learning algorithm with no eligibility trace, called TD(0), using these state representations  
720 (Fig. 3a). While TD(0) has a learning rate parameter ( $\alpha$ ), it did not influence the steady-state results,  
721 which are presented, and thus the only parameter which influenced the result was  $\gamma$ , the temporal discount  
722 factor, set to 0.925 for all simulations using a timestep of  $\Delta t = 0.2s$  (Extended Data Fig. 4 shows the  $\gamma$   
723 parameter search space). Code for generating the simulated experiments and implementing the  
724 simulations can be found at: <https://github.com/mhburrell/Qian-Burrell-2024>

#### 725 *CSC-TD model with and without ITI states*

726 We initially simulated the Conditioning, Degradation, and Cued Reward experimental conditions using  
727 the CSC-TD model, adapted from Schultz et al.<sup>25</sup>. The cue length was fixed at 1 unit of time, with time  
728 unit size set to 0.2 s, and the ISI was matched to experimental parameters at 3.5 s. Simulated cue and  
729 reward frequencies were matched to experimental parameters, separately simulating Conditioning then  
730 Degradation and Conditioning then Cued Reward. In complete serial compound, also known as tapped-  
731 delay line, each cue results in a cascade of discrete substates that completely tile the ISI. TD error and  
732 value were then modelled using a standard TD(0) implementation<sup>21</sup>, using  $\alpha = 0.1$ ,  $\gamma = 0.925$ . Reported  
733 values are the average of the last 200 instances averaged for 25 simulations. The model was run with  
734 states tiling the ISI only (CSC) or tiling the ISI and ITI until the next cue presentation (CSC with ITI  
735 states).

#### 736 *Context-TD model*

737 The Context-TD model, which is an extension of the CSC-TD model, includes context as an additional  
738 cue, but otherwise identical to the CSC simulations. For each phase (Conditioning, Degradation, Cued  
739 Reward) a separate context state was active for the entire phase, including the ISI and ITI. This  
740 corresponds to the additive Cue-Context model previously described<sup>16,17,19</sup>. TD errors reported are the  
741 average of the last 200 instances averaged for 25 simulations, except for Extinction which corresponded  
742 to third day of training.

#### 743 *Belief-State model*

744 We simulated the TD error signaling in all four conditions (Conditioning, Degradation, Cued Reward,  
745 Extinction) using a previously described belief-state TD model<sup>50</sup>. For comparison to the CSC based  
746 models described above, we had a total of 19 states, 17 capturing the ISI substates (3.5s in 0.2s  
747 increments, as in the CSC model). State 18 we termed the ‘Wait’ state and state 19 the ‘pre-transition’ or  
748 ‘pre’ state. In the Belief-State model it is assumed the animal has learned a state transition distribution.

749 We computed the transition matrix by labelling the simulated experiments with state, labelling the fixed  
750 post-US period as the Wait state and the variable ITI as the Pre state and then empirically calculating the  
751 transition matrix for that simulation. While the post-US and variable ITI periods were used to estimate the  
752 rate of transition between the Wait and Pre states, because we assumed a fixed probability of transition,  
753 these should not be considered identical – rather the implicit assumption in modeling with a fixed  
754 probability is that the time in the Wait state is a geometric random variable.

755 The belief-state model also assumes that the animal has learned a probability of distributions given the  
756 current state, encoded in an observation matrix. In our implementation there are five possible  
757 observations: Odor A, B, C, reward and null (no event). Like the transition matrix, the observation matrix  
758 was calculated empirically from the simulated experiments. Fig 3b represents schematically the state-  
759 space of the Belief-State model: Odor A (and C in Cued Reward) are observed when transitioning from  
760 Pre to the first ISI state; reward is observed in transition from the last ISI state to Wait, Odor B (and  
761 reward in Deg) are observed when transitioning from Pre to Wait. We did not consider how the details of  
762 how the transition and observation matrices may be learnt on a trial-by-trial basis as the steady-state TD  
763 errors are not dependent on this implementation. As for the other models, the TD errors reported are the  
764 average of the last 200 instances averaged over 25 simulations, except for Extinction which corresponded  
765 to the third day of training.

766 A relative value metric was used as a potential explanation of the decrease in licking during the ISI in the  
767 Degradation condition (Extended Data Fig 5). Relative value at time t was computed as value at time t (as  
768 defined and simulated by the Belief-State TD model) divided by the total value of the entire session,  
769 multiplied by the total number of rewards in a session.

770

## 771 **RNN Modeling**

772 We implemented value-RNNs, as described previously<sup>55</sup>, to model the responses in the three conditions  
773 (Conditioning, Degradation, Cued Reward). Briefly, simulated tasks were generated to match  
774 experimental parameters using a time step of 0.5s. We then trained recurrent network models, in PyTorch,  
775 to estimate value. Each value-RNN consisted of between 5 and 50 GRU cells, followed by a linear  
776 readout of value. The hidden unit activity, taken to be the RNN's state representation, can be written as  
777  $z_t = f_\phi(o_t, z_{t-1})$  given parameters  $\phi$ . The RNN's output was the value estimate  $V_t = w^\top z_t + w_0$ , for  
778  $z_t, w \in \mathbb{R}^H$  (where H is the number of hidden units) and  $V_t, w_0 \in \mathbb{R}$ . The full parameter vector  $\theta =$   
779  $[\phi \ w \ w_0]$  was learned using TD learning. This involved backpropagating the gradient of the squared error  
780 loss  $\delta_t^2 = (r_t + \gamma V_{t+1} - V_t)^2$  with respect to  $V_t$  on episodes composed of 20 concatenated trials. The  
781 timestep size was 0.5 s and  $\gamma$  was 0.83 to match the 0.925 for 0.2 s timesteps used in the TD simulations,  
782 such that both had a discount rate of 0.67 per second.

783 Prior to training, the weights and biases were initialized with the PyTorch default. To replicate the actual  
784 training process, we initially trained the RNNs on the Cond simulations, then on either the Degradation or  
785 Cued Reward conditions (Fig 6b). Training on the Cond simulations for 300 epochs on a session of  
786 10,000 trials, with a batch size of 12 episodes. Parameter updates used Adam with an initial learning rate  
787 of 0.001. To replicate the actual training process, we initially trained the RNNs on the Cond simulations,  
788 then on either the Degradation or Cued Reward conditions (Fig 6b). To simulate animals' internal timing  
789 uncertainty, the reward timing was jittered 0.5 seconds on a random selection of trials. The model

790 summary plots (Fig 6c, Extended Data Fig 6) presents mean RPE for each event. Exemplar trials shown  
791 in Fig 6 have the jitter removed for display purposes.

792 To visualize the state space used, we performed a two-step canonical correlation analysis process  
793 adapting methods used to identify long-term representation stability in the cortex<sup>74</sup>. Briefly, in each  
794 condition, we applied principal component analysis (PCA) to identify the principal components (PCs) that  
795 explained 80% of the variance (mean number of components = 4.26), then used CCA (Python package  
796 pyrcca) to project the PCs into a single space for all conditions. CCA finds linear combinations of each of  
797 the PCs that maximally correlate – allowing us to identify hidden units encoding the same information in  
798 the different RNNs. We then used the combination of PCA and CCA to create a map from hidden unit  
799 activity to a common state.

800 We measured belief  $R^2$  as previously described<sup>55</sup>. For each simulation, we calculated the beliefs from the  
801 observations of cues and rewards. We then used multivariate linear regression to decode these beliefs  
802 from hidden unit activity. To evaluate model fit, calculate the total variance explained as:  $R^2 = 1 -$   
803  $\frac{Var(B-B_{est})}{Var(B)}$ , where  $B_{est}$  is the estimate from the regression and  $Var(X) = \frac{1}{T} \sum_{t=1}^T \|x_t - \bar{x}\|^2$ .

804

## 805 ANCCR model

806 The ANCCR model is a recent alternative explanation of dopamine function<sup>31</sup>. While two previous  
807 studies have tested contingency degradation with ANCCR, they did not include the cued-reward controls.  
808 We implemented the ANCCR model using the code provided on the repository site  
809 (<https://github.com/namboodirilab/ANCCR>) and matching the simulation parameters to the experiment.  
810 We used the set of parameter values used in the previous studies, trying both Jeong et al., (2022) and Garr  
811 et al. (2023). The total parameter space searched was: T ratio = 0.2-2,  $\alpha = 0.01-0.3$ ,  $k = 0.01-1$  or  $1/(\text{mean}$   
812  $\text{inter-reward interval})$ ,  $w = 0-1$ , threshold = 0.1-0.7,  $\alpha_R = 0.1-0.3$ . The presented results use the  
813 parameters from Garr et al. (2023), as they were a better fit (T ratio = 1,  $\alpha = 0.2$ ,  $k = 0.01$ ,  $w = 0.4$ ,  
814 threshold = 0.7,  $\alpha_R = 0.1$ ). Additionally, we varied the weight of prospective and retrospective processes  
815 ( $w$ ) to examine whether the data can be explained better by choosing a specific weight. Data presented are  
816 the last 200 instances averaged for the same 25 simulations used in the TD simulations.

817

## 818 Outcome Specific Degradation Modeling

819 To model outcome-specific degradation we adapted both our Belief-State model and RNN models. For  
820 the Belief-State, we estimated the transition and observation matrix for the experiments described in Garr  
821 et al., 2023 (depicted in Extended Data Fig 8a) as described for our experiment, using a time step of 1s.  
822 As there were two rewarded trial types, we had representations of two ISI periods (termed ISI 1 and ISI 2,  
823 depicted in Extended Data Fig 8). The model was initially trained on the liquid reward (setting  $r=1$  when  
824 observing liquid reward,  $r=0$  when observing food reward) and the average TD error calculated for each  
825 trial type. We then trained on only the food reward. The total TD error was calculated as the absolute  
826 difference between the TD error on each reward type.

827 For the RNN models, we similarly adjusted the timestep to 1s and trained on simulated experiments to  
828 match the experimental parameters. Rather than training separately, the model was trained on both  
829 simultaneously, training to produce an estimate of the value of the liquid reward and an estimate of the

830 food reward, then using the 2-dimensional vector TD error to train the model. This ensures a single state  
831 space is used to solve for both reward types. Total TD error was calculated as the absolute difference on  
832 each reward type post-hoc.

833

### 834 **Statistical analysis**

835 Data analysis was performed using third party packages (e.g. Scipy, Statsmodel, etc.) in Python. All code  
836 used for analysis is available on request. Our behavioral data and dopamine response data have passed the  
837 normality test. For statistical comparisons of the mean, we used Student's *t*-test with a significance  
838 threshold of 0.05, adjusted with the Bonferroni correction. We used Welch's test for dopamine response  
839 to various events due to unequal variance between groups. Paired *t*-tests were conducted when the same  
840 mouse's performance was being compared across two different sessions. No statistical methods were used  
841 to predetermine sample sizes, but our sample sizes are similar to those reported in previous publications.  
842 The assumptions of the *t*-test were tested using the Shapiro-Wilk test to check for normality and Levene's  
843 test to check for equal variance.

844 **Material availability.**

845

846 **Data availability.** The behavioral and fluorometry data will be shared at a public deposit source.

847

848 **Code availability.** The model code will be attached as Supplementary Data. All other conventional codes  
849 used to obtain the results will be available from a public deposit source.

## 850 Reference

851

- 852 1. Rescorla, R. A. Probability of shock in the presence and absence of CS in fear conditioning. *J. Comp.*  
853 *Physiol. Psychol.* **66**, 1–5 (1968).
- 854 2. Rescorla, R. A. Conditioned inhibition of fear resulting from negative CS-US contingencies. *J.*  
855 *Comp. Physiol. Psychol.* **67**, 504–509 (1969).
- 856 3. Rescorla, R. A. Pavlovian conditioning. It's not what you think it is. *Am. Psychol.* **43**, 151–160  
857 (1988).
- 858 4. Gibbon, J., Berryman, R. & Thompson, R. L. Contingency spaces and measures in classical and  
859 instrumental conditioning. *J. Exp. Anal. Behav.* **21**, 585–605 (1974).
- 860 5. Hallam, S. C., Grahame, N. J. & Miller, R. R. Exploring the edges of Pavlovian contingency space:  
861 An assessment of contingency theory and its various metrics. *Learn. Motiv.* **23**, 225–249 (1992).
- 862 6. Cheng, P. W. From covariation to causation: A causal power theory. *Psychol. Rev.* **104**, 367 (1997).
- 863 7. Gallistel, C. R., Craig, A. R. & Shahan, T. A. Contingency, contiguity, and causality in conditioning:  
864 Applying information theory and Weber's Law to the assignment of credit problem. *Psychol. Rev.*  
865 **126**, 761–773 (2019).
- 866 8. Jenkins, H. M. & Ward, W. C. JUDGMENT OF CONTINGENCY BETWEEN RESPONSES AND  
867 OUTCOMES. *Psychol. Monogr.* **79**, SUPPL 1:1-17 (1965).
- 868 9. Allan, L. G. Human contingency judgments: rule based or associative? *Psychol. Bull.* **114**, 435–448  
869 (1993).
- 870 10. Bermudez, M. A. & Schultz, W. Responses of amygdala neurons to positive reward-predicting  
871 stimuli depend on background reward (contingency) rather than stimulus-reward pairing (contiguity).  
872 *J. Neurophysiol.* **103**, 1158–1170 (2010).
- 873 11. Griffiths, T. L. & Tenenbaum, J. B. Structure and strength in causal induction. *Cognit. Psychol.* **51**,  
874 334–384 (2005).
- 875 12. Gershman, S. J. & Ullman, T. D. Causal implicatures from correlational statements. *PloS One* **18**,  
876 e0286067 (2023).
- 877 13. Balsam, P. D. & Gallistel, C. R. Temporal maps and informativeness in associative learning. *Trends*  
878 *Neurosci.* **32**, 73–78 (2009).
- 879 14. Papini, M. R. & Bitterman, M. E. The role of contingency in classical conditioning. *Psychol. Rev.* **97**,  
880 396–403 (1990).
- 881 15. Kamin, L. Selective association and conditioning. in *Fundamental issues in associative learning* 42–  
882 64 (1969).
- 883 16. Rescorla, R. A. & Wagner, A. R. A Theory of Pavlovian Conditioning: Variations in the  
884 Effectiveness of Reinforcement and Nonreinforcement. in *Classical conditioning II: current research*  
885 *and theory* (eds. Black, A. & Prokasy, W.) 64–99 (1972).



- 886 17. Pearce, J. M. & Bouton, M. E. Theories of associative learning in animals. *Annu. Rev. Psychol.* **52**,  
887 111–139 (2001).
- 888 18. Bouton, M. E. *Learning and Behavior: A Contemporary Synthesis*. (Sinauer Associates, Inc.,  
889 Sunderland, MA, 2007).
- 890 19. Madarasz, T. J. *et al.* Evaluation of ambiguous associations in the amygdala by learning the structure  
891 of the environment. *Nat. Neurosci.* **19**, 965–972 (2016).
- 892 20. Gershman, S. J. Context-dependent learning and causal structure. *Psychon. Bull. Rev.* **24**, 557–565  
893 (2017).
- 894 21. Sutton, R. S. & Barto, A. G. *Reinforcement Learning: An Introduction*. vol. 1 (MIT Press,  
895 Cambridge, MA, 1998).
- 896 22. Mnih, V. *et al.* Human-level control through deep reinforcement learning. *Nature* **518**, 529–533  
897 (2015).
- 898 23. Sutton, R. S. Learning to predict by the methods of temporal differences. *Mach. Learn.* **3**, 9–44  
899 (1988).
- 900 24. Sutton, R. S. & Barto, A. G. Time-derivative models of Pavlovian reinforcement. in *Learning and*  
901 *computational neuroscience: Foundations of adaptive networks* (eds. Gabriel, M. & Moore, J.) 497–  
902 537 (The MIT Press, Cambridge, MA, US, 1990).
- 903 25. Schultz, W., Dayan, P. & Montague, P. R. A neural substrate of prediction and reward. *Science* **275**,  
904 1593–1599 (1997).
- 905 26. Cohen, J. Y., Haesler, S., Vong, L., Lowell, B. B. & Uchida, N. Neuron-type-specific signals for  
906 reward and punishment in the ventral tegmental area. *Nature* **482**, 85–88 (2012).
- 907 27. Watabe-Uchida, M., Eshel, N. & Uchida, N. Neural Circuitry of Reward Prediction Error. *Annu. Rev.*  
908 *Neurosci.* **40**, 373–394 (2017).
- 909 28. Kim, H. R. *et al.* A Unified Framework for Dopamine Signals across Timescales. *Cell* **183**, 1600-  
910 1616.e25 (2020).
- 911 29. Amo, R. *et al.* A gradual temporal shift of dopamine responses mirrors the progression of temporal  
912 difference error in machine learning. *Nat. Neurosci.* **25**, 1082–1092 (2022).
- 913 30. Niv, Y. Reinforcement learning in the brain. *J. Math. Psychol.* **53**, 139–154 (2009).
- 914 31. Jeong, H. *et al.* Mesolimbic dopamine release conveys causal associations. *Science* **378**, eabq6740  
915 (2022).
- 916 32. Garr, E. *et al.* Mesostriatal dopamine is sensitive to specific cue-reward contingencies.  
917 2023.06.05.543690 Preprint at <https://doi.org/10.1101/2023.06.05.543690> (2023).
- 918 33. K Nambodiri, V. M. & Stuber, G. D. The learning of prospective and retrospective cognitive maps  
919 within neural circuits. *Neuron* **109**, 3552–3575 (2021).
- 920 34. Escobar, M. & Miller, R. R. A Review of the Empirical Laws of Basic Learning in Pavlovian  
921 Conditioning. *Int. J. Comp. Psychol.* **17**, (2004).

- 922 35. Durlach, P. J. Role of signals for unconditioned stimulus absence in the sensitivity of autoshaping to  
923 contingency. *J. Exp. Psychol. Anim. Behav. Process.* **15**, 202–211 (1989).
- 924 36. Kim, C. K. *et al.* Simultaneous fast measurement of circuit dynamics at multiple sites across the  
925 mammalian brain. *Nat. Methods* (2016) doi:10.1038/nmeth.3770.
- 926 37. Bäckman, C. M. *et al.* Characterization of a mouse strain expressing Cre recombinase from the 3'  
927 untranslated region of the dopamine transporter locus. *Genes. N. Y. N 2000* **44**, 383–390 (2006).
- 928 38. Daigle, T. L. *et al.* A suite of transgenic driver and reporter mouse lines with enhanced brain cell type  
929 targeting and functionality. *Cell* **174**, 465–480.e22 (2018).
- 930 39. de Jong, J. W. *et al.* A Neural Circuit Mechanism for Encoding Aversive Stimuli in the Mesolimbic  
931 Dopamine System. *Neuron* **101**, 133–151.e7 (2019).
- 932 40. Menegas, W., Babayan, B. M., Uchida, N. & Watabe-Uchida, M. Opposite initialization to novel  
933 cues in dopamine signaling in ventral and posterior striatum in mice. *eLife* **6**, (2017).
- 934 41. Akam, T., Costa, R. & Dayan, P. Simple Plans or Sophisticated Habits? State, Transition and  
935 Learning Interactions in the Two-Step Task. *PLoS Comput. Biol.* **11**, e1004648 (2015).
- 936 42. Takahashi, Y. K. *et al.* Expectancy-related changes in firing of dopamine neurons depend on  
937 orbitofrontal cortex. *Nat. Neurosci.* **14**, 1590–1597 (2011).
- 938 43. Starkweather, C. K. & Uchida, N. Dopamine signals as temporal difference errors: recent advances.  
939 *Curr. Opin. Neurobiol.* **67**, 95–105 (2020).
- 940 44. Takahashi, Y. K. *et al.* Dopaminergic prediction errors in the ventral tegmental area reflect a  
941 multithreaded predictive model. *Nat. Neurosci.* **26**, 830–839 (2023).
- 942 45. Daw, N. D., Courville, A. C. & Touretzky, D. S. Representation and Timing in Theories of the  
943 Dopamine System. *Neural Comput.* **18**, 1637–1677 (2006).
- 944 46. Kobayashi, S. & Schultz, W. Influence of reward delays on responses of dopamine neurons. *J.*  
945 *Neurosci. Off. J. Soc. Neurosci.* **28**, 7837–7846 (2008).
- 946 47. Fiorillo, C. D., Newsome, W. T. & Schultz, W. The temporal precision of reward prediction in  
947 dopamine neurons. *Nat. Neurosci.* **11**, 966–973 (2008).
- 948 48. Enomoto, K. *et al.* Dopamine neurons learn to encode the long-term value of multiple future rewards.  
949 *Proc. Natl. Acad. Sci. U. S. A.* **108**, 15462–15467 (2011).
- 950 49. Masset, P. *et al.* Multi-timescale reinforcement learning in the brain. *BioRxiv Prepr. Serv. Biol.*  
951 2023.11.12.566754 (2023) doi:10.1101/2023.11.12.566754.
- 952 50. Starkweather, C. K., Babayan, B. M., Uchida, N. & Gershman, S. J. Dopamine reward prediction  
953 errors reflect hidden-state inference across time. *Nat. Neurosci.* **20**, 581–589 (2017).
- 954 51. Starkweather, C. K., Gershman, S. J. & Uchida, N. The Medial Prefrontal Cortex Shapes Dopamine  
955 Reward Prediction Errors under State Uncertainty. *Neuron* **98**, 616–629.e6 (2018).
- 956 52. Namboodiri, V. M. K. *et al.* Single-cell activity tracking reveals that orbitofrontal neurons acquire  
957 and maintain a long-term memory to guide behavioral adaptation. *Nat. Neurosci.* **22**, 1110–1121  
958 (2019).

- 959 53. Schultz, W. Dopamine reward prediction-error signalling: a two-component response. *Nat. Rev.*  
960 *Neurosci.* **17**, 183–195 (2016).
- 961 54. Lak, A., Nomoto, K., Keramati, M., Sakagami, M. & Kepecs, A. Midbrain Dopamine Neurons Signal  
962 Belief in Choice Accuracy during a Perceptual Decision. *Curr. Biol. CB* **27**, 821–832 (2017).
- 963 55. Hennig, J. A. *et al.* Emergence of belief-like representations through reinforcement learning. *PLoS*  
964 *Comput. Biol.* **19**, e1011067 (2023).
- 965 56. Bach, F. R. & Jordan, M. I. Kernel independent component analysis. *J. Mach. Learn. Res.* **3**, 1–48  
966 (2003).
- 967 57. Sussillo, D., Churchland, M. M., Kaufman, M. T. & Shenoy, K. V. A neural network that finds a  
968 naturalistic solution for the production of muscle activity. *Nat. Neurosci.* **18**, 1025–1033 (2015).
- 969 58. Gershman, S. J. A Unifying Probabilistic View of Associative Learning. *PLoS Comput. Biol.* **11**,  
970 e1004567 (2015).
- 971 59. Namboodiri, V. M. K. How do real animals account for the passage of time during associative  
972 learning? *Behav. Neurosci.* **136**, 383–391 (2022).
- 973 60. Naneix, F., Marchand, A. R., Di Scala, G., Pape, J.-R. & Coutureau, E. A role for medial prefrontal  
974 dopaminergic innervation in instrumental conditioning. *J. Neurosci. Off. J. Soc. Neurosci.* **29**, 6599–  
975 6606 (2009).
- 976 61. Piquet, R., Faugère, A. & Parkes, S. L. A hippocampo-cortical pathway detects changes in the  
977 validity of an action as a predictor of reward. *Curr. Biol.* **0**, (2023).
- 978 62. Delamater, A. R., Siegel, D. B. & Tu, N. C. Learning about reward identities and time. *Behav.*  
979 *Processes* **207**, 104859 (2023).
- 980 63. Grove, J. C. R. *et al.* Dopamine subsystems that track internal states. *Nature* **608**, 374–380 (2022).
- 981 64. Willmore, L. *et al.* Overlapping representations of food and social stimuli in mouse VTA dopamine  
982 neurons. *Neuron* **111**, 3541-3553.e8 (2023).
- 983 65. Millidge, B., Song, Y., Lak, A., Walton, M. E. & Bogacz, R. Reward-Bases: Dopaminergic  
984 Mechanisms for Adaptive Acquisition of Multiple Reward Types. 2023.05.09.540067 Preprint at  
985 <https://doi.org/10.1101/2023.05.09.540067> (2023).
- 986 66. Pearl, J. *Causality*. (Cambridge university press, 2009).
- 987 67. Baird, L. C. Advantage Updating. Technical report WL-TR-93-1146. Wright-Patterson Air Force  
988 Base. (1993).
- 989 68. Dayan, P. & Balleine, B. W. Reward, motivation, and reinforcement learning. *Neuron* **36**, 285–298  
990 (2002).
- 991 69. Schulman, J., Moritz, P., Levine, S., Jordan, M. & Abbeel, P. High-Dimensional Continuous Control  
992 Using Generalized Advantage Estimation. Preprint at <https://doi.org/10.48550/arXiv.1506.02438>  
993 (2018).

- 994 70. Pan, H.-R., Gürtler, N., Neitz, A. & Schölkopf, B. Direct Advantage Estimation. in *Advances in*  
995 *Neural Information Processing Systems* (eds. Koyejo, S. et al.) vol. 35 11869–11880 (Curran  
996 Associates, Inc., 2022).
- 997 71. Veitch, V., Wang, Y. & Blei, D. M. Using Embeddings to Correct for Unobserved Confounding in  
998 Networks. Preprint at <https://doi.org/10.48550/arXiv.1902.04114> (2019).
- 999 72. Dezfouli, A. & Balleine, B. W. Habits, action sequences and reinforcement learning. *Eur. J.*  
1000 *Neurosci.* **35**, 1036–1051 (2012).
- 1001 73. Langdon, A. J., Sharpe, M. J., Schoenbaum, G. & Niv, Y. Model-based predictions for dopamine.  
1002 *Curr. Opin. Neurobiol.* **49**, 1–7 (2018).
- 1003 74. Gallego, J. A., Perich, M. G., Chowdhury, R. H., Solla, S. A. & Miller, L. E. Long-term stability of  
1004 cortical population dynamics underlying consistent behavior. *Nat. Neurosci.* **23**, 260–270 (2020).
- 1005 75. Fehr, E. & Goette, L. Do Workers Work More if Wages Are High? Evidence from a Randomized  
1006 Field Experiment. *Am. Econ. Rev.* **97**, 298–317 (2007).
- 1007 76. Herrnstein, R. J. Relative and absolute strength of response as a function of frequency of  
1008 reinforcement. *J. Exp. Anal. Behav.* **4**, 267–272 (1961).
- 1009

## 1010 **Acknowledgements**

1011 We thank Hao Wu and Nune Martiros for technical assistance on the behavioral code design; Mitsuko  
1012 Uchida for discussion and advice on task design; Catherine Dulac, Florian Engert, and all lab members  
1013 from Naoshige Uchida lab and Venkatesh Murthy lab for discussion. This work was supported by grants  
1014 from the National Institute of Health (U19 NS113201, R01DC017311), the Simons Collaboration on  
1015 Global Brain, the Air Force Office of Scientific Research (FA9550-20-1-0413), the Human Frontier  
1016 Science Program (LT000801/2018 to S.M.), the Harvard Brain Science Initiative, and the Brain and  
1017 Behavior Research Foundation (NARSAD Young Investigator no. 30035 to S.M.).

1018

1019

## 1020 **Author information**

1021 These authors contributed equally: Lechen Qian, Mark Burrell

1022

## 1023 **Contributions**

1024 L.Q., N.U., and V.N.M conceived the conceptual framework and designed the behavioral tasks and  
1025 recording experiments. L.Q. conducted all experiments and data analysis. S.M. established the multifiber  
1026 photometry system and supplied the transgenic mice. M.B., N.U. and L.Q. discussed the modeling  
1027 framework. M.B. constructed all the TD learning models and conducted the analysis. J.H. constructed  
1028 RNN models. The RNN modeling results were analyzed by M.B., J.H. and L.Q. The results were  
1029 discussed and interpreted by L.Q., N.U., M.B., J.H, S.G. and V.N.M. The manuscript was written by  
1030 L.Q., M.B., and N.U. and all other authors provided feedback.

1031

## 1032 **Corresponding authors**

1033 Correspondence to Naoshige Uchida ([uchida@mcb.harvard.edu](mailto:uchida@mcb.harvard.edu))

1034

1035

## 1036 **Competing interests**

1037 The authors declare no competing financial interests.

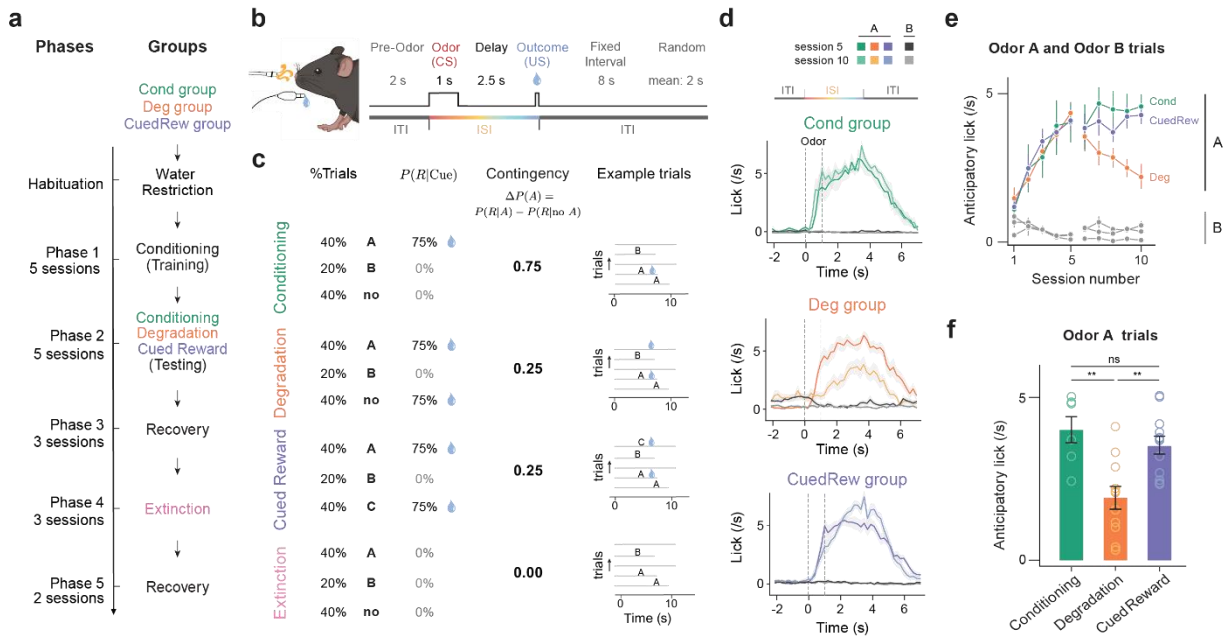
1038

1039

1040

1041

1042



1043

1044 **Figure 1 | Dynamic changes in lick response to olfactory cues across different phases of Pavlovian contingency**  
1045 **learning task.**

1046 (a) Experimental design. Three groups of mice subjected to four unique conditions of contingency learning. All  
1047 animals underwent Phases 1 and 2. Deg group additionally underwent Phases 3-5.

1048 (b) Trial timing.

1049 (c) Trial parameters per condition. In Conditioning, Degradation and Cued Reward, Odor A predicts 75% chance of  
1050 reward (9  $\mu$ L water) delivery, Odor B indicates no reward. In Degradation, blank trials were replaced with  
1051 uncued rewards (75% reward probability). In Cued Reward, these additional rewards were cued by Odor C. In  
1052 Extinction, no rewards were delivered.

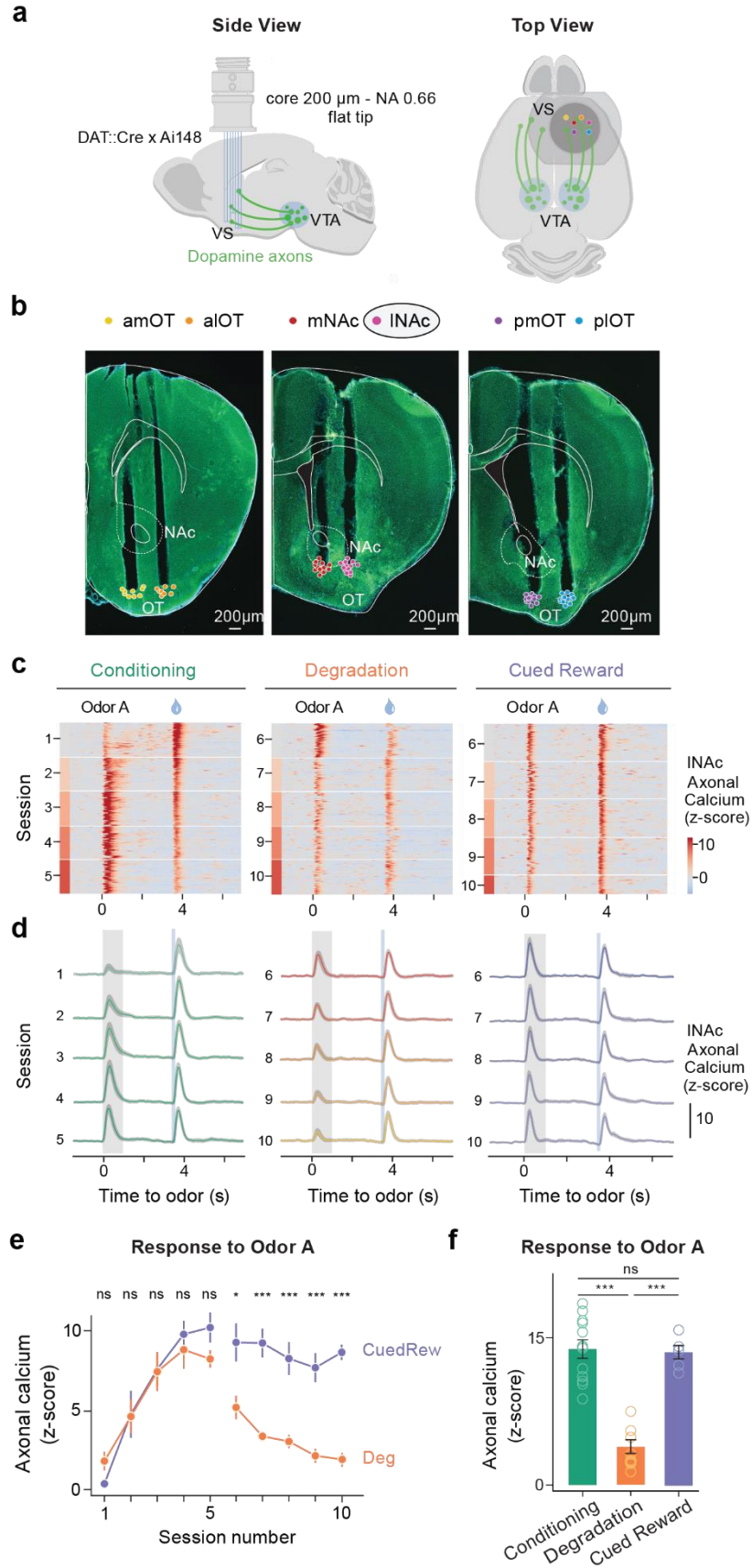
1053 (d) PSTH of average licking response of mice in three groups to the onset of Odor A and Odor B from the last  
1054 session of Phase 1 (session 5) and Phase 2 (session 10). Shaded area is standard error of the mean (SEM).  
1055 Notably, the decreased licking response during ISI and increased during ITI in Deg group. (green, Cond group,  
1056  $n = 6$ ; orange, Deg group,  $n = 11$ ; purple, CuedRew group,  $n = 12$  mice).

1057 (e) Average lick rate in 3s post-cue (Odor A or B) by session. Error bars represent SEM.

1058 (f) Average lick rate in 3s post Odor A in final session of each condition. Asterisks denote statistical significance:  
1059 ns,  $P > 0.05$ ; \*\*,  $P < 0.01$ , Student's  $t$ -test, indicating a significant change in licking behavior to Odor A in Deg  
1060 group across sessions.

1061

1062



1064

1065 **Figure 2 | Dopamine axonal activity recordings show different responses to rewarding cues in Degradation**  
1066 **and Cued Reward conditions**

1067 (a) Configuration of multifiber photometry recordings. Coronal section from one DAT::cre x Ai148 mouse showing  
1068 tracts for multiple fibers in the VS. Data recorded from INAc is used in the following analysis. INAc, Lateral  
1069 nucleus accumbens; mNAc, Medial NAc; a lot, anterior lateral olfactory tubercle; plot, posterior lateral OT;  
1070 amOT, anterior medial OT; pmOT, posterior medial OT.

1071 (b) Heatmap from two example mice (mouse 1, left two panels, mouse 2, right panel) illustrating the z-scored  
1072 dopamine axonal signals in Odor A rewarded trials (rows), aligned to the onset of Odor A for three conditions.

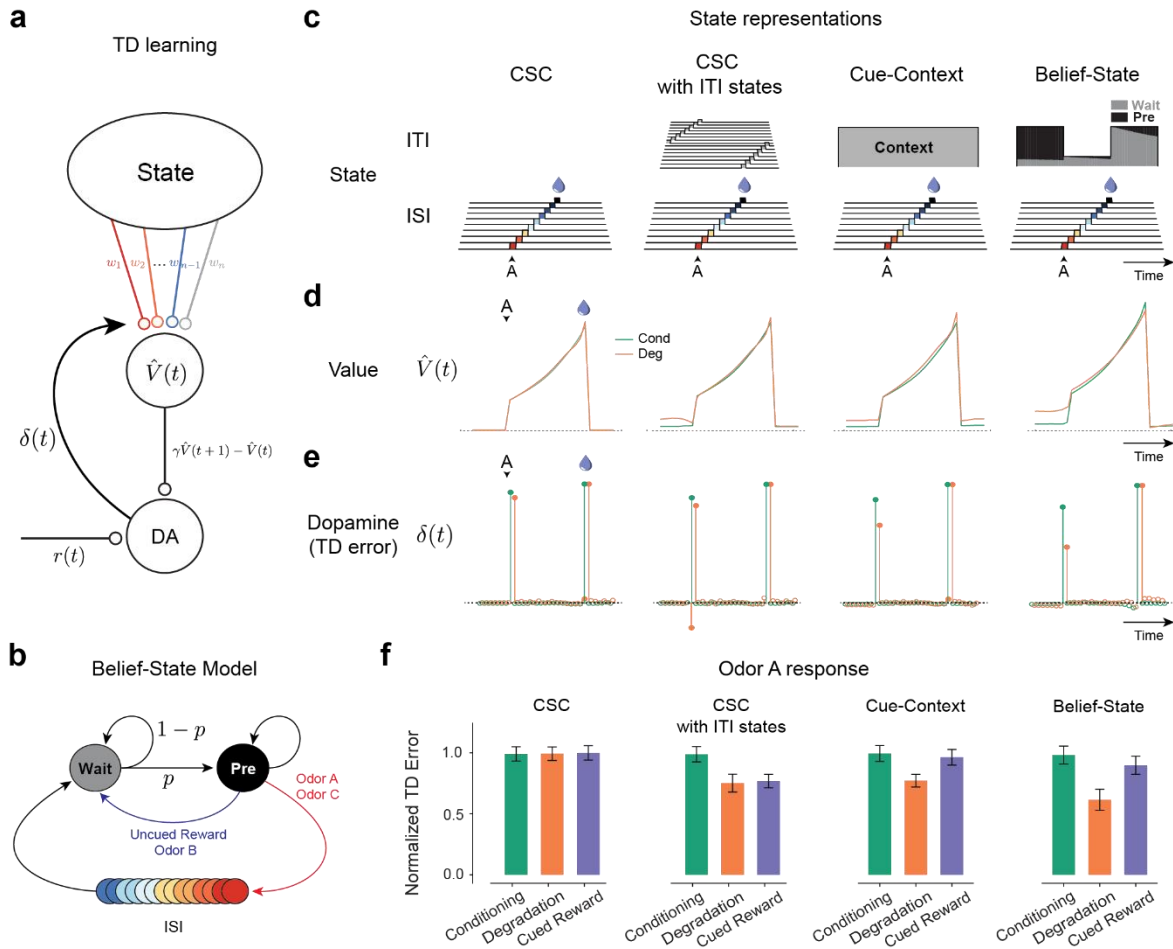
1073 (c) Population average z-scored dopamine axonal signals in response to Odor A and water delivery. Shaded areas  
1074 represent SEM.

1075 (d) Mean peak dopamine axonal signal (z-scored) of Odor A response by sessions for the Deg group (orange) and  
1076 the CuedRew group (purple). Error bars are SEM. \*,  $P < 0.05$ ; \*\*\*,  $P < 0.001$ , Welch's *t*-test.

1077 (e) Mean peak dopamine axonal signal (z-scored) for the last session in Phase 1 (Conditioning) and 2 (Degradation  
1078 and Cued Reward) for both Deg and CuedRew groups. Error bars represent SEM. ns,  $P > 0.05$ ; \*\*\*,  $P < 0.001$ ,  
1079 Welch's *t*-test.

1080





1081

1082 **Figure 3 | TD learning models can explain dopamine responses in contingency degradation with appropriate**  
 1083 **ITI representation.**

1084 (a) Temporal Difference Zero, TD(0), model – The state representation determines value. The difference in value  
 1085 between the current and gamma-discounted future state plus the reward determines the reward prediction error  
 1086 or dopamine. This error drives updates in the weights.

1087 (b) Belief-State Model: After the ISI, the animal is in the Wait state, transitioning to the pre-transition (‘Pre’) state  
 1088 with fixed probability  $p$ . Animal only leaves Pre state following the observation of odor or reward.

1089 (c) State representations: from the left, Complete Serial Compound (CSC) with no ITI representation, CSC with ITI  
 1090 states, Cue-Context model and the Belief-State model.

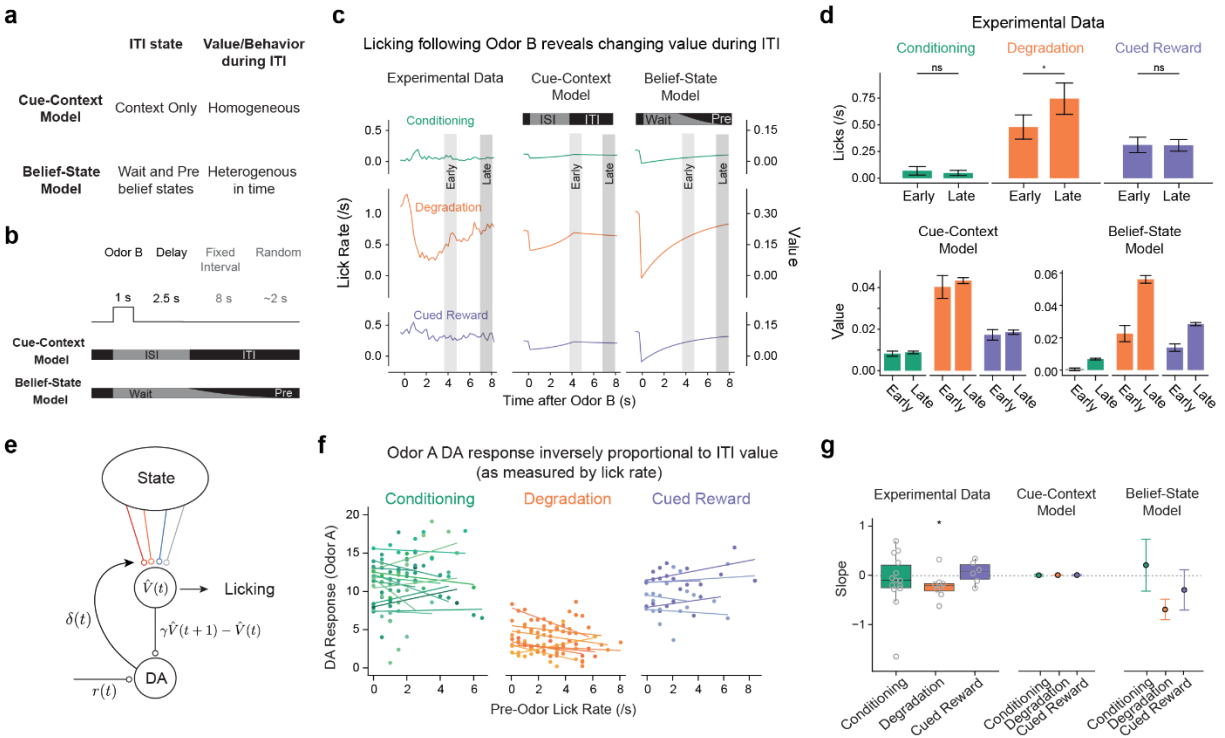
1091 (d) Value in Odor A trials of each state representation using TD(0) for Conditioning and Degradation conditions

1092 (e) TD error is the difference in value plus the reward.

1093 (f) Mean normalized TD error of Odor A response from 25 simulated experiments. Error bars are SD.

1094

1095



1096

1097

1098

**Figure 4 | Belief-State model, but not Cue-Context model, explains variance in behavior and dopamine responses.**

1099

(a) Cue-Context model and Belief-State model differ in their representation of the ITI.

1100

(b) Odor B predicts no reward and at least 10 s before the start of the next trial.

1101

(c) Odor B induces a reduction in licking, particularly in the Degradation condition, which matches the pattern of

1102

value in the Belief-State model better than the Cue-Context model.

1103

(d) Quantified licks (top) from experimental data in early (3.5-5s) and late (7-8s) post cue period. Error bars are

1104

SEM, \*,  $P < 0.05$ , paired  $t$ -test. Value from Cue-Context and Belief-State model for the same time period, error

1105

bars are SD.

1106

(e) If licking is taken as a readout of value, then ITI licking should be inversely correlated with dopamine.

1107

(f) Per animal linear regression of Odor A dopamine response (z-score axonal calcium) on lick rate in 2s before cue

1108

delivery.

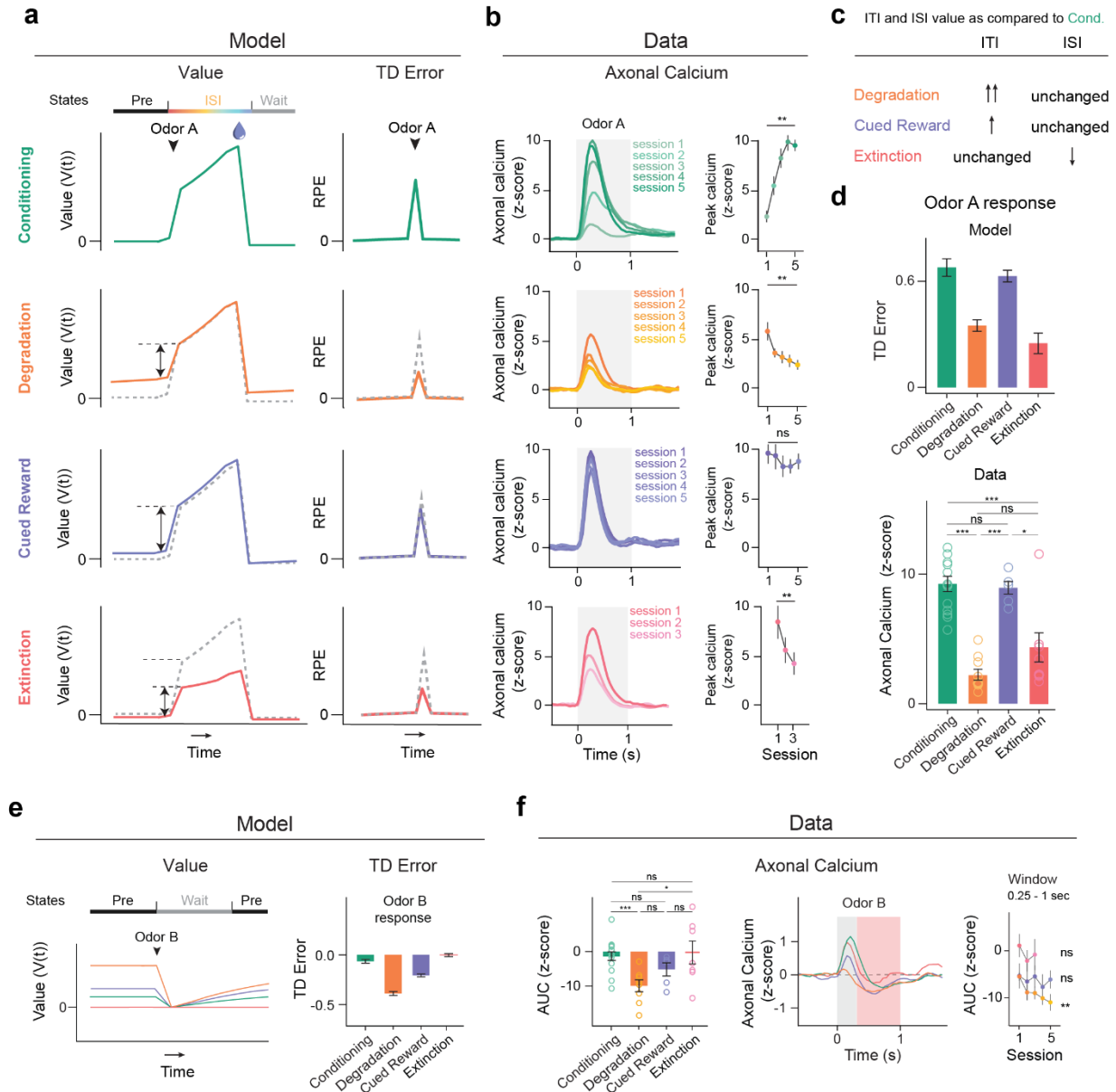
1109

(g) Summarized slope coefficients from experimental data (left) and models (right). Boxplot shows median and

1110

IQR, one sample  $t$ -test.

1111



1112

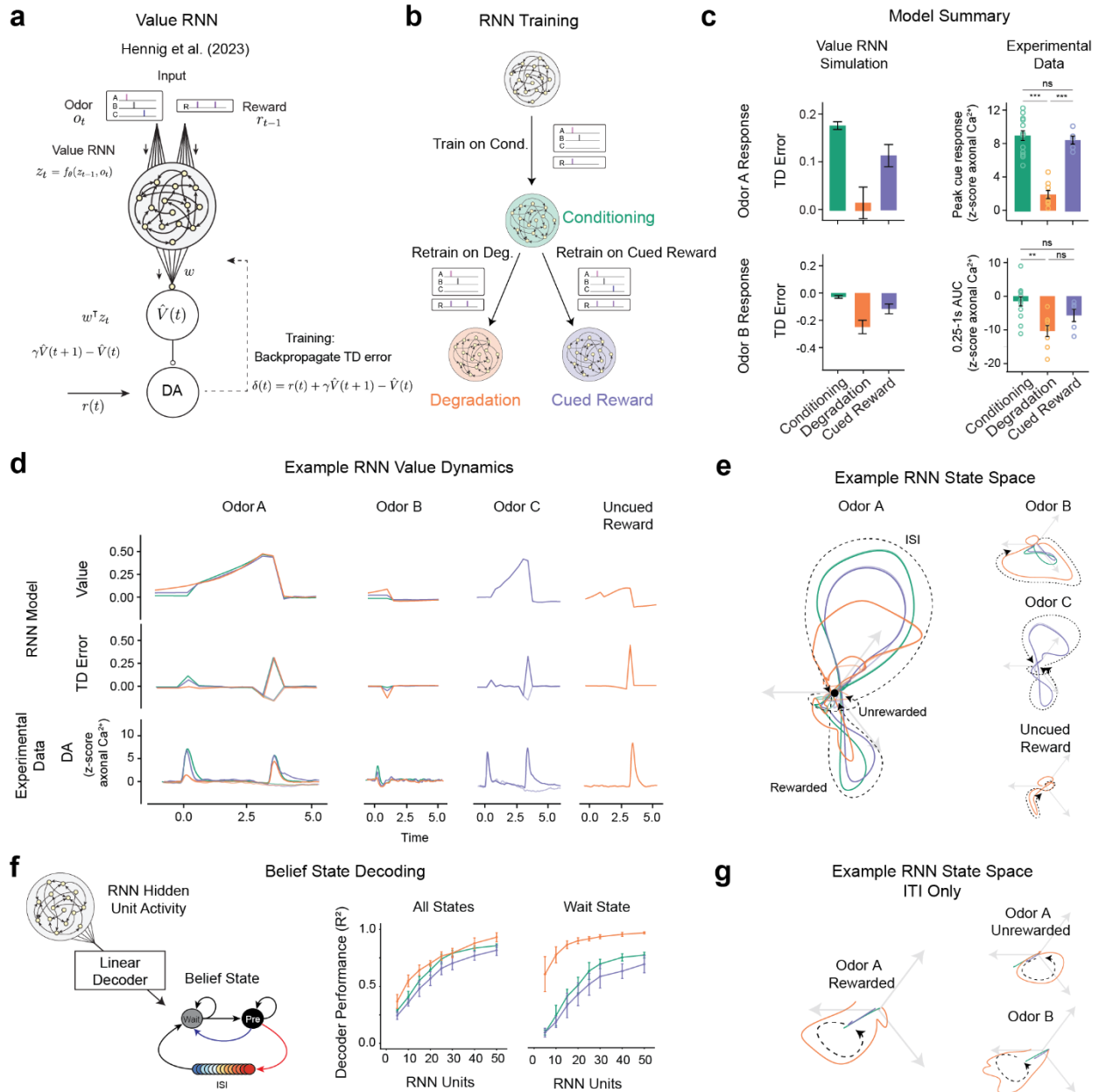
1113 **Figure 5 | Belief-State model's predictions recapitulate additional experimental data**

1114 (a) Plots averaged from one representative simulation of Odor A rewarded trial ( $n = 4,000$  simulated trials) for four  
 1115 distinct conditions using the Belief-State model. Graphs are for the corresponding value function (left) and TD  
 1116 error (right) of cue response for Odor A rewarded trials.

1117 (b) Signals from dopamine axons (mean) across multiple sessions of each condition (left). Mean peak dopamine  
 1118 axonal calcium signal (z-scored) for the first to last session in Phase 2 for four contingency conditions (right).  
 1119 Error bars represent SEM. ns,  $P > 0.05$ ; \*\*,  $P < 0.01$ , Student's paired  $t$ -test. The Belief-State model captures the  
 1120 modulation of Odor A dopamine response in all conditions.

1121 (c) Degradation, Cued Reward and Extinction conditions differ in how their ITI and ISI values change compared to  
 1122 Conditioning phase.

- 1123 (d) Mean peak TD error by Belief-State model and dopamine axonal signal (z-scored) to Odor A for four distinct  
1124 conditions. Error bars represent SEM. ns,  $P > 0.05$ ; \*,  $P < 0.05$ ; \*\*\*,  $P < 0.001$ , Welch's  $t$ -test. The model's  
1125 prediction captured well the pattern in the dopamine data.
- 1126 (e) Averaged traces from a representative simulation of Odor B trial ( $n = 4,000$  simulated trials) across four distinct  
1127 conditions using the Belief-State model. Graphs are for the value function and TD errors of cue response for  
1128 Odor B trials.
- 1129 (f) Z-scored dopamine axonal signals to Odor B quantified from the red shaded area to quantify the later response  
1130 only. Bar graph (left) shows mean z-scored Odor B AUC from 0.25s-1s response from the last session of each  
1131 condition. Error bars are SEM. \*  $P < 0.05$ ; \*\*\*,  $P < 0.001$ , Welch's  $t$ -test. Line graph (right) shows mean z-  
1132 scored AUC over multiple sessions for each condition. Statistical analysis was performed on data from the first  
1133 and last sessions of these conditions. Error bars are SEM.
- 1134
- 1135
- 1136
- 1137
- 1138
- 1139
- 1140



1141

1142

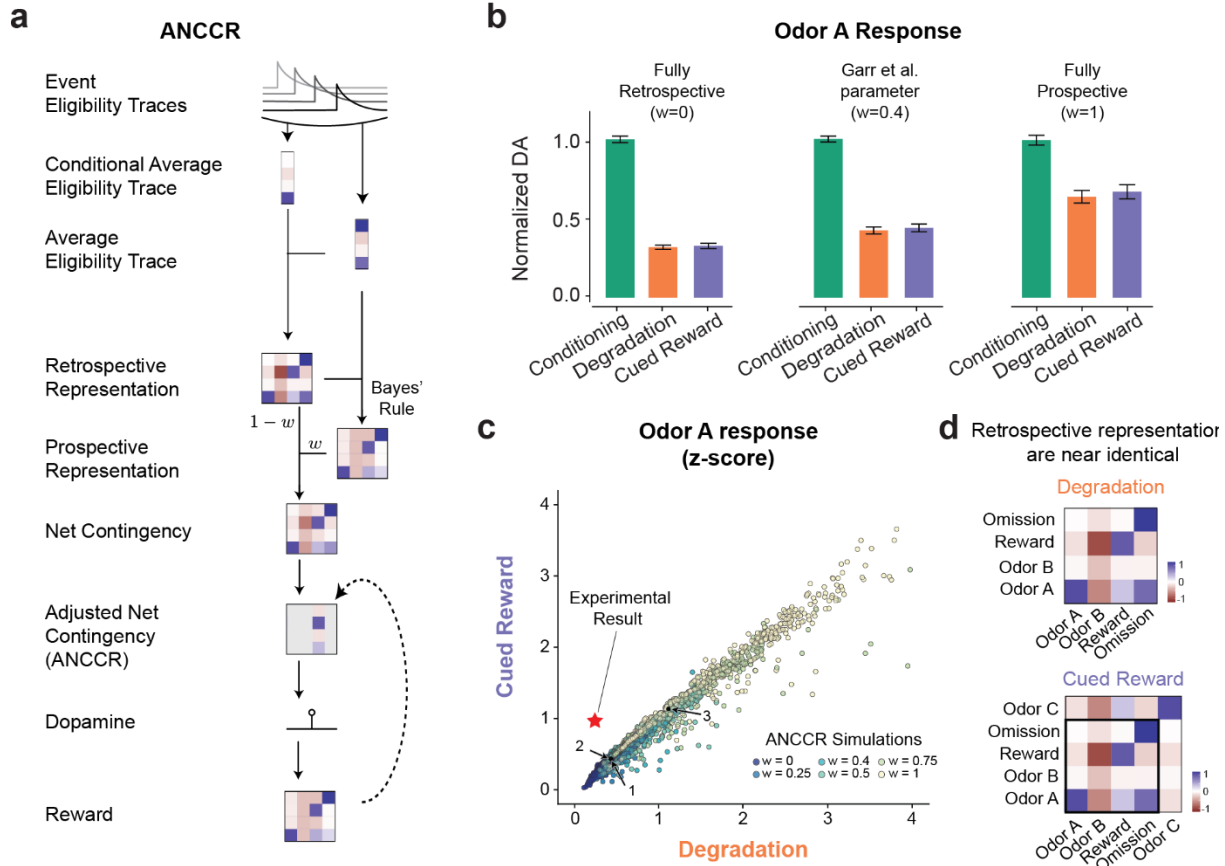
1143 **Figure 6 | Value-RNNs recapitulate experimental results using state-spaces akin to hand-crafted Belief-State**  
 1144 **model**

1145 (a) The Value-RNN replaces the hand-crafted state space representation with an RNN that is trained only on the  
 1146 observations of cues and rewards. The TD error is used to train the network.

1147 (b) RNNs were initially trained on simulated Conditioning experiments, before being retrained on either  
 1148 Degradation or Cued Reward conditions.

1149 (c) The predictions of the RNN models (mean, error bars: SD) closely match the experimental results.

- 1150 (d) Example value, TD error, and corresponding average experimental data from a single RNN simulation. Notably,  
1151 decreased Odor A response is explained by increased value in the pre-cue period.
- 1152 (e) Hidden neuron activity projected into 3D space using CCA from the same RNNs used in (d). The Odor A ISI  
1153 representation is similar in each of the three conditions, and similar to the Odor C representation. Odor B  
1154 representation is significantly changed in the Degradation condition.
- 1155 (f) Correspondence between RNN state space and Belief-State model. A linear decoder was trained to predict  
1156 beliefs using RNN hidden unit activity. With increasing hidden layer size, the RNN becomes increasingly  
1157 belief-like. The improved performance of the decoder for the Degradation condition is explained by better  
1158 decoding of the Wait state. Better Wait state decoding is explained by altered ITI representation:
- 1159 (g) Same RNNs as in (d) and (e), hidden unit activity projected into state-space as (e) for the ITI period only  
1160 reveals ITI representation is significantly different in the Degradation case.
- 1161
- 1162



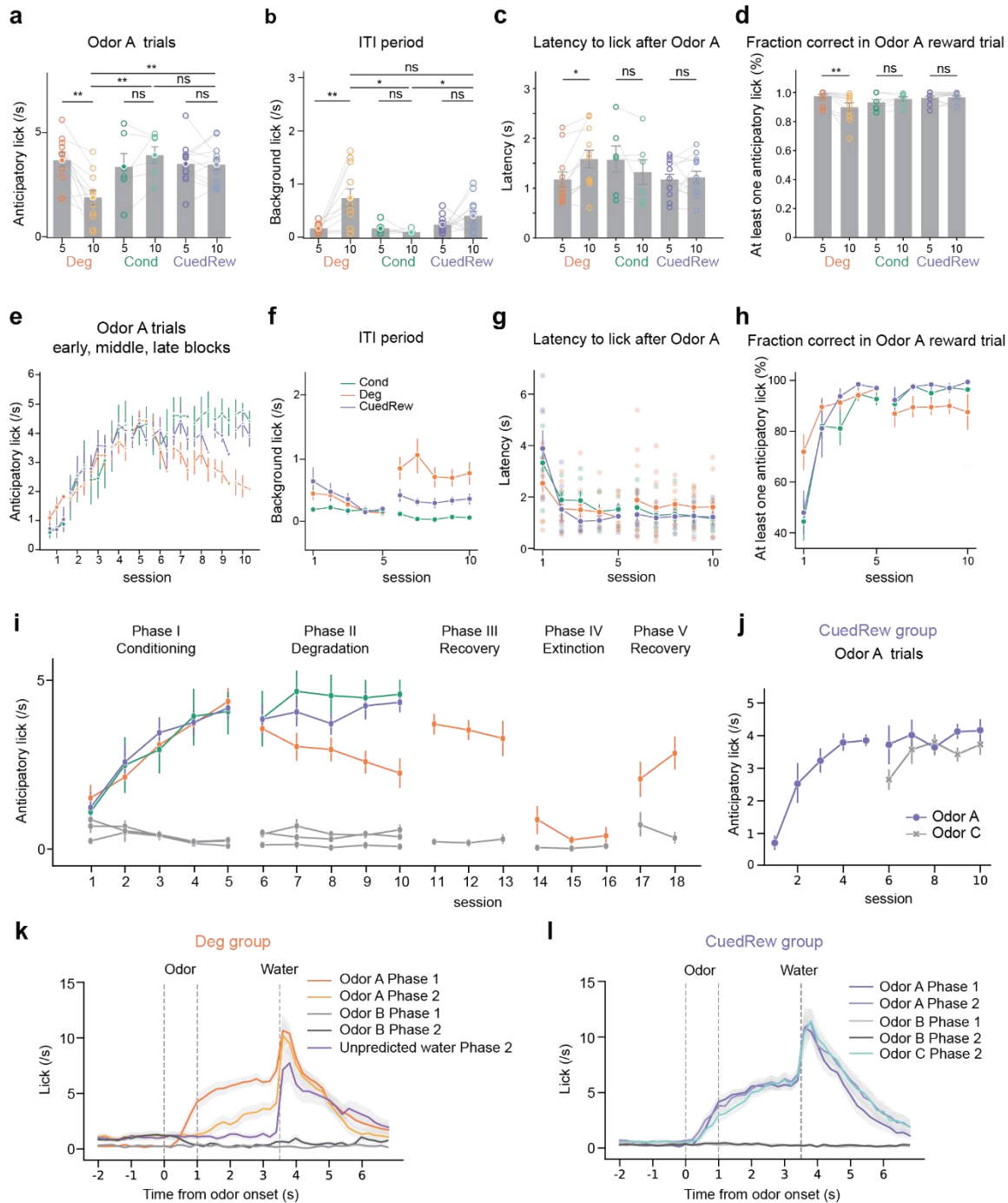
1163

1164 **Figure 7 | ANCCR does not explain the experimental results:**

- 1165 (a) Simplified representation of ANCCR model. Notably the first step is to estimate retrospective contingency
- 1166 using eligibility traces.
- 1167 (b) Simulations of the same virtual experiments used in Figure 3 using ANCCR, using the parameters in Garr et al.,
- 1168 2023 varying the prospective-retrospective weighting parameter ( $w$ ). Error bars are SD. In all cases the
- 1169 predicted odor A response is similar in the Degradation and Cued Reward conditions.
- 1170 (c) No parameter combination explains the experimental result. Searching 21,000 parameter combinations across
- 1171 six parameters ( $T$  ratio = 0.2-2,  $\alpha$  = 0.01-0.3,  $k$  = 0.01-1 or  $1/(\text{mean inter-reward interval})$ ,  $w$  = 0-1,
- 1172 threshold = 0.1-0.7,  $\alpha R$  = 0.1-0.3). Experimental result plotted as a star. Previously used parameters (Garr et
- 1173 al., 2023 as 1, Jeong et al., 2022 as 2 and 3) indicated. Dots are colored by the prospective-retrospective
- 1174 weighting parameter ( $w$ ), which has a strong effect on the magnitude of Phase 2 response relative to Phase 1.
- 1175 (d) As the contingency is calculated as the first step, and the contingencies are similar in Degradation and Cued
- 1176 Reward conditions, there is little difference in the retrospective contingency representation between the two
- 1177 conditions, explaining why regardless of parameter choice ANCCR predicts similar responses.

1178

## 1179 Extended Data Figures



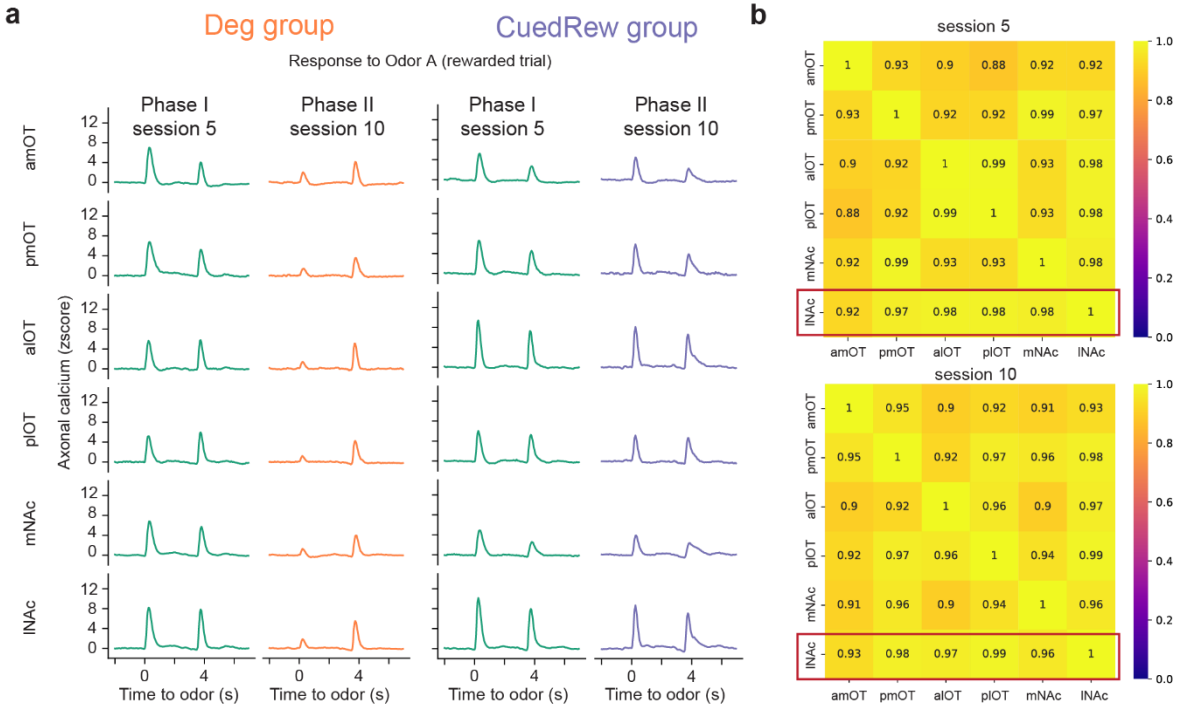
1180

### 1181 Extended Data Fig. 1 | Population Average Behavior per session

1182 (a, b, c, d) Bar graphs comparing the average number of licks to Odor A during the first 3s post-stimulus (a) and during  
 1183 ITI (b), latency to lick (c), and fraction correct (d) in the final sessions of phase 1 and phase 2 for Deg, Cond, and  
 1184 CuedRew groups. Error bars represent SEM. Asterisks denote statistical significance: ns  $p > 0.05$ , \*\* $p < 0.01$ ,  
 1185 paired Student's t-test

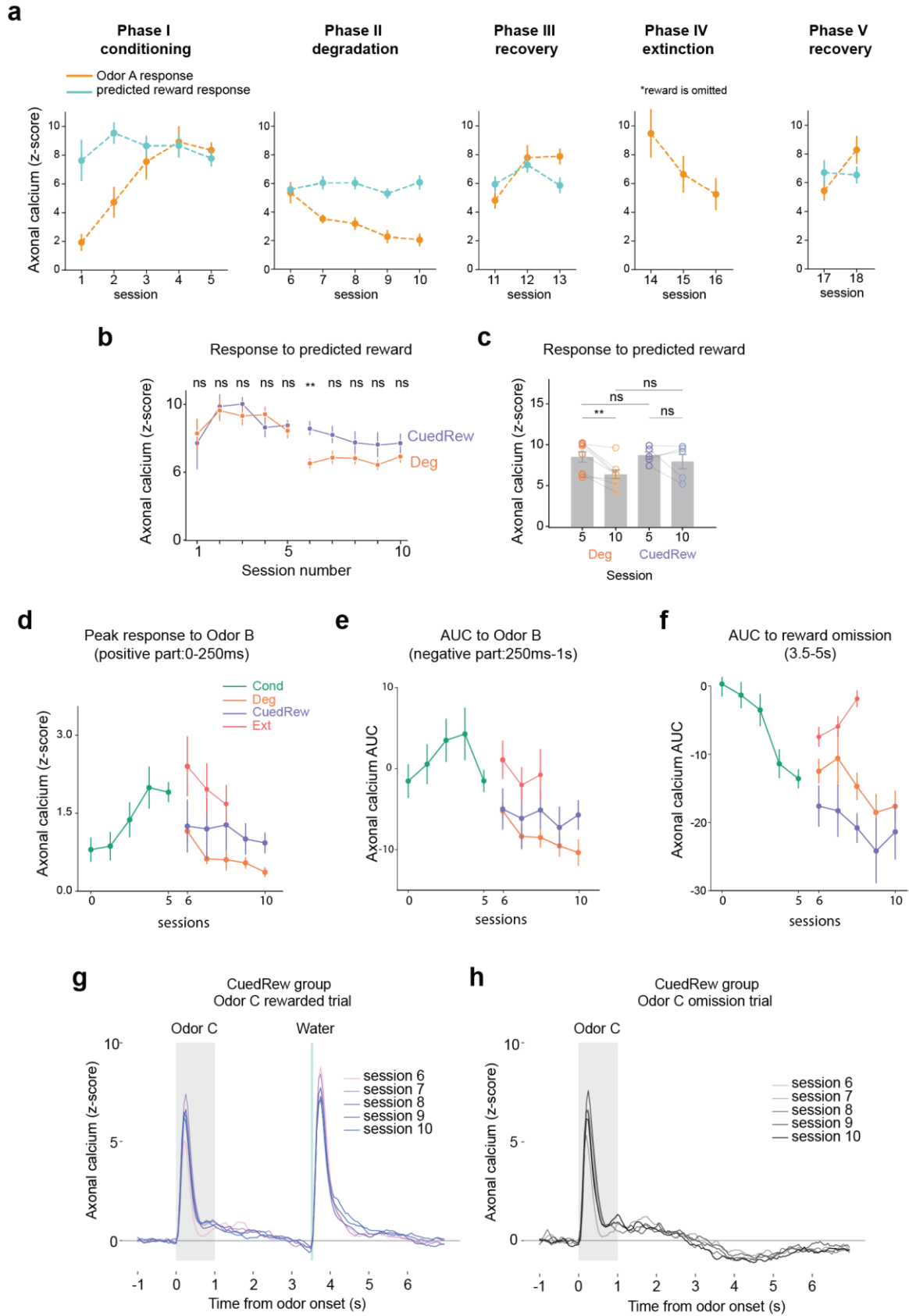


- 1186 (e) Session-wise variation in anticipatory licking for Odor A trials, broken down into early, middle, and late blocks,  
1187 for all groups.
- 1188 (f, g, h). Line graphs showing the average number of licks to Odor A (colored) during ITI (g), latency to lick after  
1189 Odor A and fraction correct in Odor A trials for each session in the Conditioning, Degradation, and Cued Reward  
1190 phase (Deg group – orange, n = 11; Cond group – green, n = 6; CuedRew – purple, n=12 mice).
- 1191 (i) Anticipatory licking rate in Odor A trials (colored) and in Odor B trials (grey) across multiple phases: Conditioning  
1192 (Phase I), Degradation (Phase II), Recovery (Phase III), Extinction (Phase IV), and post-Extinction Recovery  
1193 (Phase V).
- 1194 (j) Anticipatory licking to Odor C develops quickly compared to Odor A, potentially reflecting generalization.
- 1195 (k, l) PSTH showing the average licking response of mice in Deg group (k) and CuedRew group (l) to the various  
1196 events. The response is time-locked to the odor presentation (time 0). The shaded area indicates the standard error  
1197 of the mean (SEM).



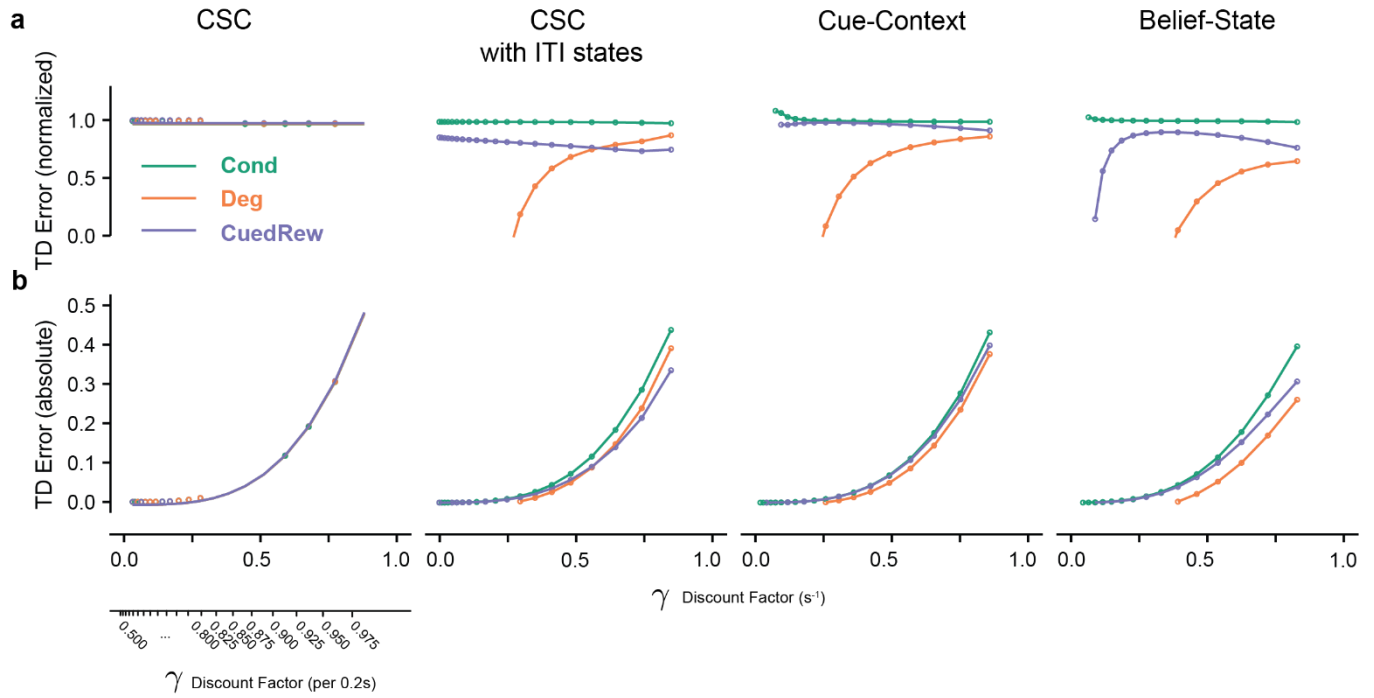
1199 **Extended Data Fig. 2 | Dopamine responses are highly correlated across recording sites**

- 1200 (a) Averaged dopamine axonal responses to Odor A during rewarded trials for both Deg group and CuedRew group,  
1201 depicted for Phase I session 5 and Phase II session 10 across all recorded sites.
- 1202 (b) Correlation matrix for averaged dopamine responses to Odor A during rewarded trials, comparing across sites  
1203 from the Deg groups during sessions 5 and 10. Cosine similarity was calculated by averaging z-scored  
1204 responses across trials within animals, then across animals and then computing the cosine similarity between  
1205 each recording site.
- 1206 (c) Population average dopamine responses to Odor A in rewarded trials across sessions 1 to 10 for both Deg and  
1207 CuedRew groups, detailing the changes in response through Phase I and Phase II.  
1208



1210 **Extended Data Fig. 3 | Population Average Dopamine Response per session**

- 1211 (a) Mean peak dopamine axonal signal (z-scored) of cue response (orange) and reward response (cyan) in Odor A  
1212 rewarded trial by sessions for the Deg group across multiple phases: Conditioning (Phase I), Degradation (Phase  
1213 II), Recovery (Phase III), Extinction (Phase IV), and post-Extinction Recovery (Phase V). Error bars are SEM.
- 1214 (b) Mean peak dopamine axonal signal (z-scored) of reward response in Odor A trials by sessions for the Deg group  
1215 (orange) and the CuedRew group (purple). Error bars are SEM. ns  $P < 0.05$ ,  $**P < 0.001$ , Student's  $t$ -test .
- 1216 (c) Mean peak dopamine axonal signal (z-scored) for the last session in Phase 1 and 2 for both Deg and CuedRew  
1217 groups. Error bars represent SEM. ns,  $P > 0.05$ ;  $***$ ,  $P < 0.001$ , paired  $t$ -test.
- 1218 (d, e, f) Mean peak dopamine axonal signal (z-scored) across sessions for four distinct conditions, represented for  
1219 various events.
- 1220 (g) Response to Odor C (rewarded) and (h) Odor C (omission), population average per session
- 1221



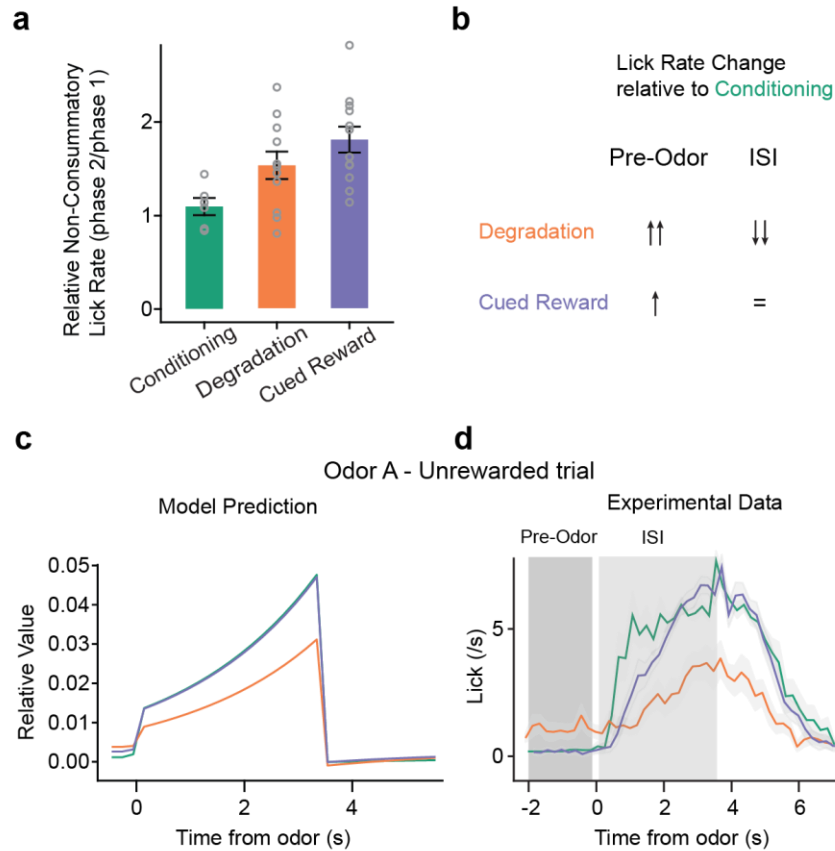
1222

1223

1224 **Extended Data Fig. 4 | Discount Factor determines modeled contingency degradation effect size**

1225 Influence of discount factor ( $\gamma$ ) on relative predicted odor A response relative to Conditioning (a) or absolute (b),  
1226 where reward size = 1 for four models presented in Figure 3. Bottom right scale showing discount factor converted to  
1227 step size (0.2s), other axes use per second discount. Tested range: 0.5-0.975 discount per 0.2s in 0.025 steps.

1228



1229

1230 **Extended Data Fig. 5 | Relative value explains decreased anticipatory licking during ISI during contingency**  
 1231 **degradation**

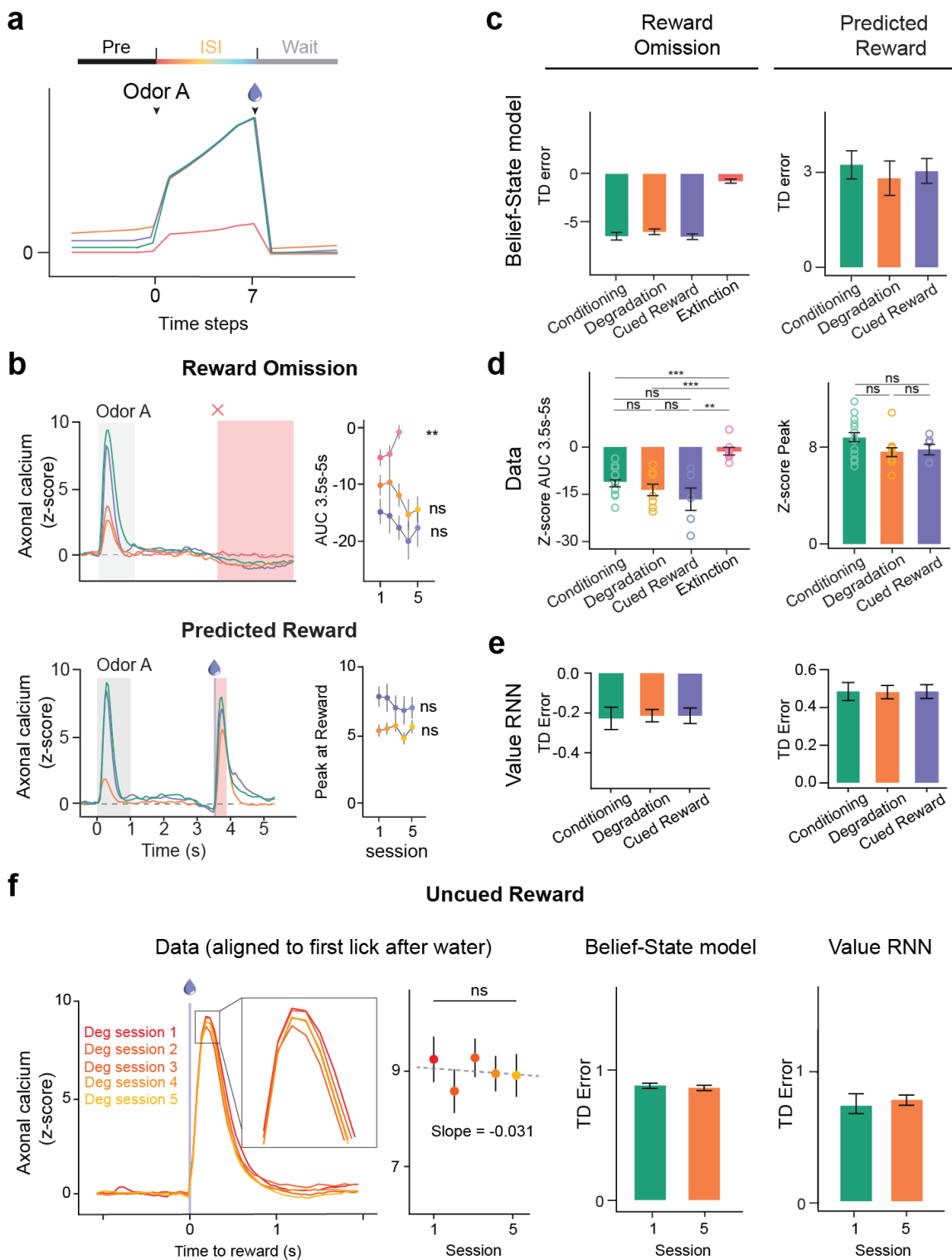
1232

1233 (a) If each lick carries a small, fixed effort cost, a rational agent will lick proportionally to the total amount  
 1234 of rewards<sup>75,76</sup>. Plot show mean non-consummatory lick rate normalized to the Conditioning phase, suggesting  
 1235 that the Degradation and Cued Reward conditions elicit approximately twice the lick rate of the Conditioning  
 1236 condition, and thus proportional to the total reward quantity. Consummatory licks were considered any licks  
 1237 occurring in the 2 seconds following reward delivery.

1238 (b) Summary of lick rate changes relative to the Conditioning phase during the pre-odor period and the inter-stimulus  
 1239 interval (ISI).

1240 (c) Average relative value (current value/session total value, scaled by total reward) during odor A trial derived from  
 1241 the Belief-State model. Relative value, which is increased in the pre-odor period and thus decreased during the  
 1242 ISI, accounts for the change in licking pattern during unrewarded (and thus without consummatory licks) odor A  
 1243 trials.

1244 (d) Experimental data showing the actual lick rates recorded during Odor A unrewarded trials, compared over time,  
 1245 which aligns with the assumptions and predictions made in a,b, and c.



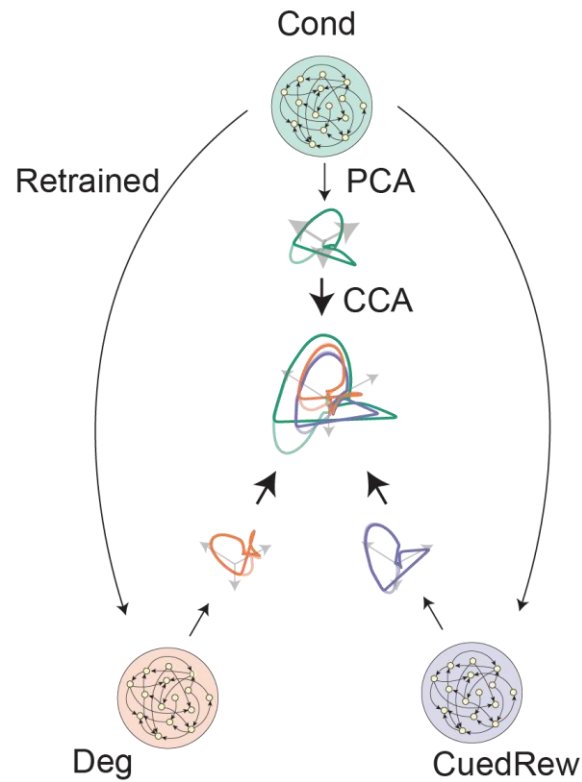
1246

1247 **Extended Data Fig. 6 | Comparison of reward and omission responses between experimental data, Belief-State**  
 1248 **model and value-RNN predictions**



- 1249 (a) Plots averaged from one representative simulation of Odor A rewarded trial ( $n = 4,000$  simulated trials) for four  
1250 distinct conditions using the Belief-State model. Graphs are for the corresponding value function of Odor A  
1251 rewarded trials, with Pre state, ISI state and Wait state annotated.
- 1252 (b) Z-scored DA axonal signals to reward omission and predicted reward following Odor A quantified from the red  
1253 shaded area. Line graphs (right) shows mean z-scored response over multiple sessions for each condition.  
1254 Statistical analysis was performed on data from the first and last session of these conditions. Error bars are  
1255 SEM. ns,  $P > 0.05$ ; \*\*,  $P < 0.01$ , paired  $t$ -test.
- 1256 (c) The predictions of the Belief-State model for reward omission and predicted reward (mean, error bars: SD).
- 1257 (d) The experimental data for reward omission and predicted reward (mean, error bars: SEM). ns,  $P > 0.05$ ; \*\*,  $P <$   
1258  $0.01$ ; \*\*\*,  $P < 0.001$ , Welch's  $t$ -test.
- 1259 (e) The predictions of the Value-RNN models for reward omission and predicted reward (mean, error bars: SD).
- 1260 (f) The experimental data, TD error prediction by Belief-State model and Value-RNN model for uncued reward  
1261 response in Degradation condition. While the Belief-State model captured the downward trend in response  
1262 magnitude, none of the three statistical tests showed significant changes.
- 1263
- 1264

1265

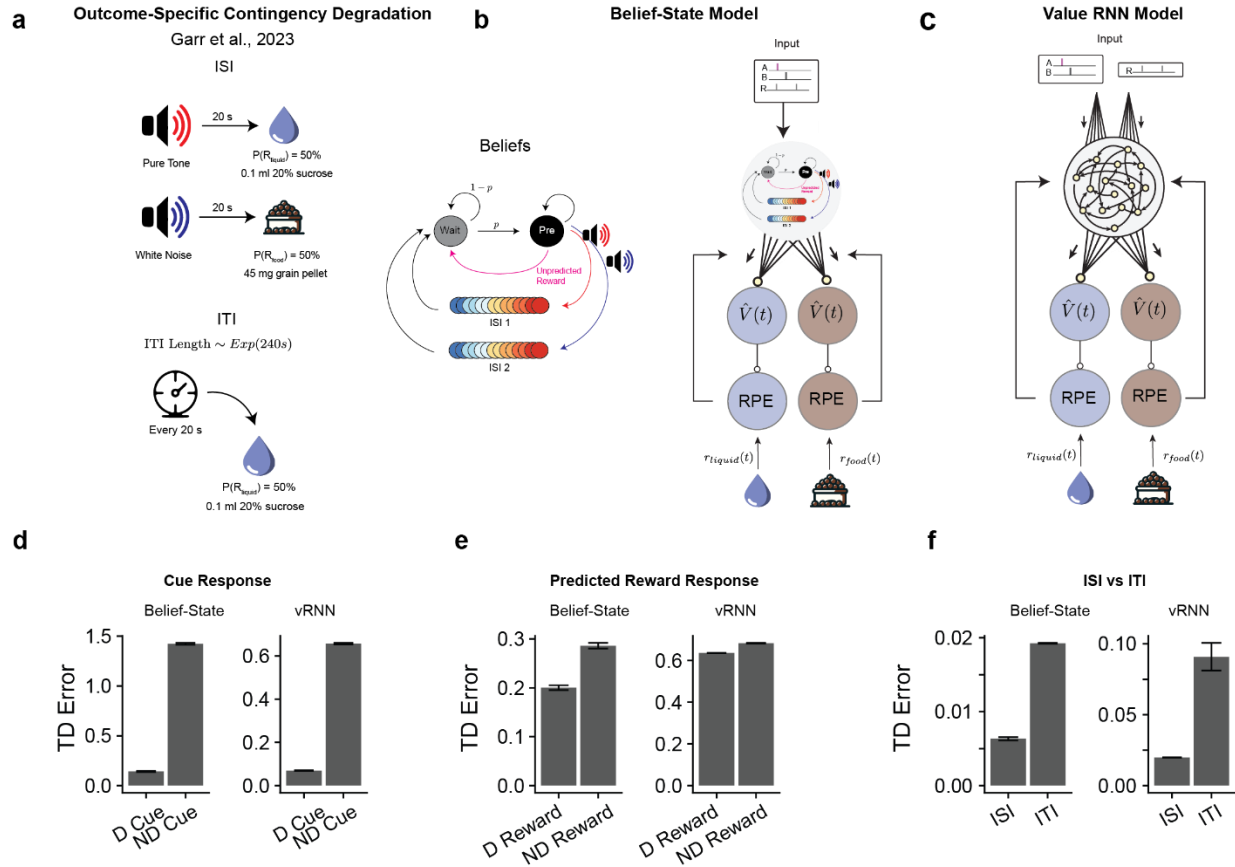


1266

1267 **Extended Data Fig. 7 | Methodology for visualizing state space from hidden unit activity**

1268 Illustration for visualizing common state space of RNN models. RNN hidden unit activity was first projected into  
1269 principal component space, then canonical correlation analysis was used to align between different conditions.

1270



1271

1272 **Extended Data Fig. 8 | Outcome-specific contingency degradation explained by Belief-State model and Value-**  
1273 **RNN model.**

1274 (a) Experimental design of Garr et al., two cues predicted either a liquid or food reward. During degradation, every  
1275 20 s the liquid reward was delivered with 50% probability. The ITI length was drawn from an exponential  
1276 distribution with mean of 4 minutes.

1277 (b) Belief-State model design. The Belief-State model was extended to include a second series of ISI substates to  
1278 reflect the two types of rewarded trials. The model was then independently trained on the liquid reward and  
1279 food reward.

1280 (c) The value-RNN model design – as (b) but replacing the Belief-State model with the value-RNN, using a vector-  
1281 valued RPE as feedback, with each channel reflecting one of the reward types.

1282 (d-f) Summary of predicted RPE responses from Belief-State Model and Value-RNN (vRNN). The RPE was  
1283 calculated as the absolute difference between the liquid RPE and food RPE. Other readout functions (e.g. weighted  
1284 sum) produce similar results. Both model predictions match experimental results with degraded (D) cue (panel d)  
1285 and degraded reward (e) having a reduced dopamine response versus non-degraded (ND). Furthermore, average  
1286 RPE during ISI (3 seconds after cue on) and ITI (3 seconds before ITI) capture measured experimental trend.  
1287 Error bars are SEM.

1288

1289

**1290 Extended Data Video 1: State Space trajectories**

1291 Animation of trajectories in CCA space from RNN presented Figure 6e. In sequence, trajectories showing Odor A  
1292 (rewarded), Odor A (unrewarded), Odor A (Rewarded and Unrewarded), Odor B and then all at once for the three  
1293 conditions. Real time speed multiple indicated top right. ITI length is extended from training/actual experiment to  
1294 demonstrate the return to original ('Pre') state in Conditioning and Cued Reward but the delayed return in Degradation  
1295 condition.