

# 1 **Revealing abrupt transitions from goal-directed to habitual behavior**

2

3 Sharlen Moore<sup>1,#</sup>, Zyan Wang<sup>1,#</sup>, Ziyi Zhu<sup>1,2,3</sup>, Ruolan Sun<sup>4</sup>, Angel Lee<sup>1</sup>, Adam Charles<sup>3,4</sup>,

4

Kishore V. Kuchibhotla<sup>1,2,3,4,\*</sup>

5

6 1 Department of Psychological and Brain Sciences, Krieger School of Arts and Sciences, Johns Hopkins  
7 University, Baltimore, MD, USA.

8 2 The Solomon H. Snyder Department of Neuroscience, Johns Hopkins School of Medicine, Baltimore,  
9 Maryland, USA.

10 3 Johns Hopkins Kavli Neuroscience Discovery Institute, Johns Hopkins University, Baltimore, MD, USA.

11 4 Department of Biomedical Engineering, Whiting School of Engineering, Johns Hopkins University,  
12 Baltimore, Maryland, USA.

13

14

15 #Equal contribution

16 \*Correspondence: [kkuchib1@jhu.edu](mailto:kkuchib1@jhu.edu)

17

18 Keywords: goal-directed, habit, motivation, learning

19

## 20 **Abstract**

21 A fundamental tenet of animal behavior is that decision-making involves multiple  
22 'controllers.' Initially, behavior is goal-directed, driven by desired outcomes, shifting later  
23 to habitual control, where cues trigger actions independent of motivational state. Clark  
24 Hull's question from 1943 still resonates today: "Is this transition abrupt, or is it gradual  
25 and progressive?"<sup>1</sup> Despite a century-long belief in gradual transitions, this question  
26 remains unanswered<sup>2,3</sup> as current methods cannot disambiguate goal-directed versus  
27 habitual control in real-time. Here, we introduce a novel 'volitional engagement' approach,  
28 motivating animals by palatability rather than biological need. Offering less palatable  
29 water in the home cage<sup>4,5</sup> reduced motivation to 'work' for plain water in an auditory  
30 discrimination task when compared to water-restricted animals. Using quantitative  
31 behavior and computational modeling<sup>6</sup>, we found that palatability-driven animals learned  
32 to discriminate as quickly as water-restricted animals but exhibited state-like fluctuations  
33 when responding to the reward-predicting cue—reflecting goal-directed behavior. These

34 fluctuations spontaneously and abruptly ceased after thousands of trials, with animals  
35 now always responding to the reward-predicting cue. In line with habitual control, post-  
36 transition behavior displayed motor automaticity, decreased error sensitivity (assessed  
37 via pupillary responses), and insensitivity to outcome devaluation. Bilateral lesions of the  
38 habit-related dorsolateral striatum<sup>7</sup> blocked transitions to habitual behavior. Thus,  
39 'volitional engagement' reveals spontaneous and abrupt transitions from goal-directed to  
40 habitual behavior, suggesting the involvement of a higher-level process that arbitrates  
41 between the two.

42

## 43 **Main text**

44 Humans and other animals are often thought to be creatures of habit. When driving, for  
45 example, we are initially told that the color of a traffic light should guide our actions: green  
46 to 'go' and red to 'stop'. Through practice, we learn purposefully and are driven by the  
47 conscious goal to follow the rules of the road. Over time, these rules become automatic;  
48 without deliberation, we will push the gas pedal on a green light and the brake pedal for  
49 a red light. This transition from goal-directed to habitual control has long been assumed  
50 to be gradual<sup>8-14</sup>. More specifically, an initial goal-directed action (R) in response to a cue  
51 (S) yields a desired outcome (O) which then slowly evolves into a habit where the cue  
52 elicits the action (S-R) without necessarily having the goal in mind<sup>15,16</sup>. The automatization  
53 of decisions can be thought of as an efficient way to offload well-learned contingencies to  
54 free up resources for more flexible, goal-directed learning<sup>17</sup>. The formation and  
55 perseverance of habits, however, can also be maladaptive with neural circuits being co-  
56 opted in substance use disorders or compulsive behaviors<sup>18-20</sup>. Understanding the exact  
57 time course of habit formation is critical to disentangling its neural basis and could help  
58 inform future interventional strategies for combating habit-related disorders.

59         The assumption of a slow, gradual shift from goal-directed to habitual control  
60 underpins current models of learning and informs most approaches to understanding the  
61 neurobiological basis of habit formation. To date, however, the nature, timing, and  
62 properties of the transition between controllers have been challenging to pinpoint due to  
63 methodological constraints. In rodents, the gold standard for assessing whether a

64 behavior is under goal-directed or habitual control at a specific time point exploits the  
65 observation that goal-directed actions are sensitive to the outcome<sup>21,22</sup> (e.g., animals will  
66 only perform an action when the reward is desired) while habitual behavior is less  
67 sensitive to the outcome (e.g., animals will continue to perform said action even if the  
68 reward is not explicitly desired). This behavioral characterization relies on defining  
69 habitual behavior as the loss or absence of goal-directed control<sup>23,24</sup>. Sensitivity to the  
70 outcome has been successfully operationalized in laboratory testing with 'outcome  
71 devaluation'<sup>18,25</sup> procedures in which a reward is devalued (through satiety or taste  
72 aversion). Outcome devaluation, or a related alternative called contingency degradation,  
73 is typically implemented at set time points outside of the normal training regimen (e.g.,  
74 the middle and end of a multi-day training period). To date, no approach exists to  
75 disambiguate between goal-directed and habitual control in real-time, during training<sup>26,27</sup>.  
76 A complementary approach in the study of habit formation in rodents is to exploit distinct  
77 reinforcement schedules to bias goal-directed, or habitual behavior<sup>28</sup>. Animals under  
78 distinct reinforcement schedules are then tested for habitual behavior with outcome  
79 devaluation. The use of outcome devaluation, contingency degradation, and distinct  
80 reinforcement schedules, though powerful, inherently limit assessing the nature, timing,  
81 and properties of the transition between goal-directed and habitual control in individual  
82 animals due to the discrete test sessions and cohort-level comparisons. Can we  
83 behaviorally identify habitual transitions in real-time and during training? Is the transition  
84 slow or sudden? What are the characteristics of these transitions in individual animals?  
85 Addressing these questions requires a new behavioral approach that assesses the  
86 decision mode *en passant*, without discrete 'test' sessions, without biasing behavior to  
87 one or the other process, and without impacting the ongoing learning process.

88 Here, we present such an approach. We reasoned that animals are usually highly  
89 motivated to perform tasks because water and/or food are restricted and only made  
90 available during the task. This leads to a sustained ceiling motivation driven by the  
91 animals' need to obtain their daily food or water intake in a short period of time. In such  
92 situations, animals remain highly engaged throughout a task irrespective of the underlying  
93 decision mode. We hypothesized that if animals were motivated mainly by a taste  
94 preference, rather than a biological need, we could track naturalistic fluctuations in their

95 motivation for the preferred reward by fostering variability in reward-seeking. Under goal-  
96 directed control, task engagement would wax and wane naturalistically ('volitional  
97 engagement'), due to palatability-driven motivation, which would lead to reduced  
98 responding to the reward-predicting cue; in contrast, under habitual control (in which  
99 animals are less sensitive to the outcome), the S-R nature of the behavior would drive  
100 high and stable responding to the reward-predicting cue despite ongoing changes in the  
101 underlying desire for the outcome.

102

### 103 **A palatability-based approach reduces motivation to consume water**

104 In the home cage, we gave mice *ad libitum* access to water laced with citric acid (CA)  
105 (**Fig. 1ai**, CA yellow) which makes it slightly acidic to the taste but still fulfills hydration  
106 needs<sup>4,5</sup>. Before instrumental training, CA mice lost significantly less weight than mice on  
107 a standard water restriction (WR85) protocol (**Fig. 1a<sub>ii</sub>**) (WR85 = 17.8%±2.3%, CA =  
108 8.9%±1.9% average and std weight loss respectively, p=0.000055, Wilcoxon rank sum  
109 test). This difference was maintained throughout training (**Extended Data Fig. 1a**)  
110 (Wilcoxon rank sum test, p=0.000055), while no differences were observed in the animals'  
111 initial weight (**Extended Data Fig. 1b**) (Wilcoxon rank sum test, p=0.97). Before mice  
112 begin discriminative auditory training (**Fig. 1ai**), they first experience 2 days of  
113 instrumental training in which they learned to make an instrumental response (lick) to  
114 subsequently receive a small reward ('lick training', 3 µl plain water droplet) (**Fig. 1ai**,  
115 lavender). During this initial session, we assessed licking patterns as a proxy of  
116 motivation. CA mice executed fewer licks per session (**Fig. 1b-c**, and **Extended Data**  
117 **Fig. 1c**, yellow) (Wilcoxon rank sum test, p=0.021), exhibited strikingly different lick  
118 patterns (**Fig. 1c**), and obtained significantly fewer rewards compared to WR85 (**Fig. 1d**)  
119 (Wilcoxon rank sum test, p=0.00079) while maintaining the same lick frequency when  
120 engaged in licking (**Extended Data Fig. 1d**) (Wilcoxon rank sum test, p=0.67). This  
121 suggests that under the CA protocol, mice exhibit reduced motivation for plain water.

122

123

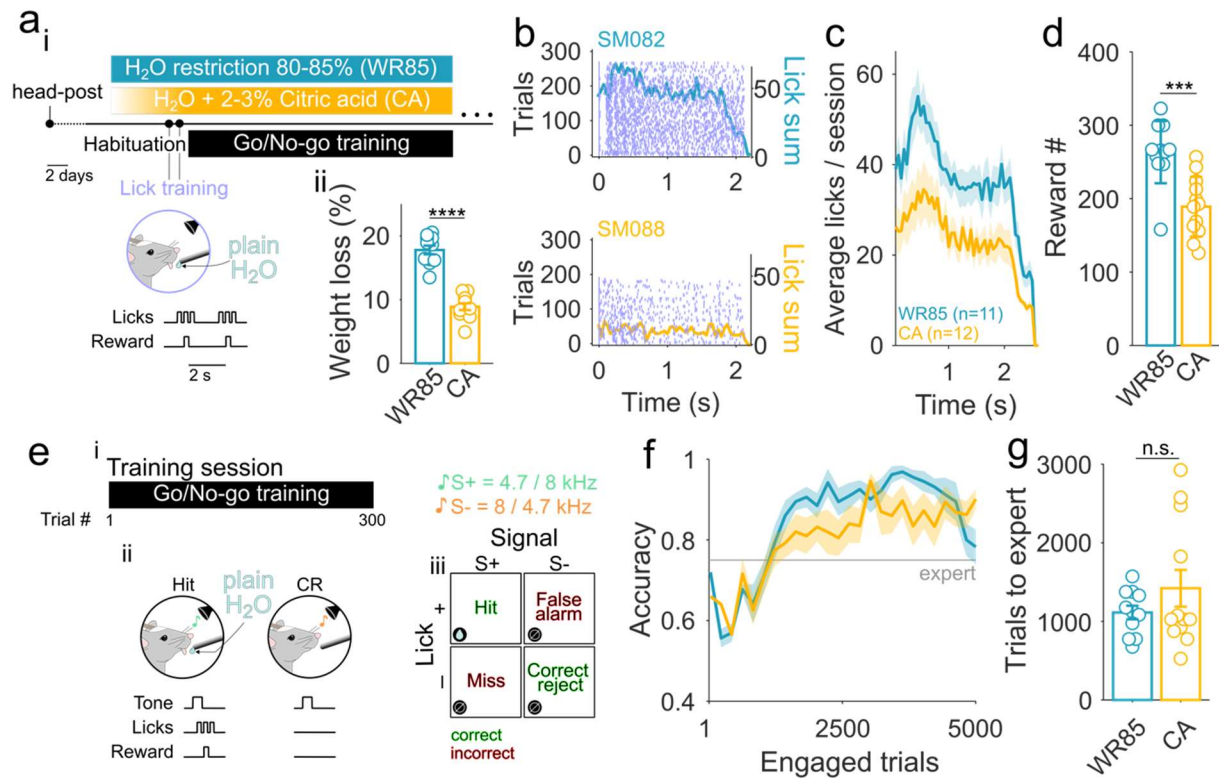
124

125

126

127

128



**Figure 1. Palatability-based motivation reduces water consumption without impacting learning trajectories.** **a<sub>i</sub>**, Protocol outline: after head-post implantation, mice underwent a habituation period and introduction to the water restriction paradigms. ‘Control’ mice underwent a common (maintenance at 80-85% original body weight) water restriction paradigm (WR85, blue). CA mice had *ad libitum* access to a water bottle with a low percentage of a less palatable hydration source (citric acid laced water) in their home cage (CA, yellow) to which they were progressively introduced (starting with 0.5%, reaching maximum 3%). After a few days, both cohorts underwent two days of lick training (lavender), followed by an auditory Go/No-Go training (black). **a<sub>ii</sub>**, CA mice lost significantly less weight compared to WR85 (Wilcoxon rank-sum test,  $p=0.000055$ , WR85  $n=11$ , CA  $n=12$ ). **b**, Lick raster plots during lick training for a WR85 (upper) and a CA (lower) exemplar animal. The right y-axis (color lines) represents a PSTH-like sum of licks for each animal. **c**, Average PSTH-like licking in one session, showing that CA mice lick less ( $n=11$  WR85, and  $n=12$  CA mice). **d**, CA mice obtain significantly fewer rewards compared to WR85 (Wilcoxon rank sum test,  $p=0.00079$ ) in one lick training session ( $n=11$  WR85, and  $n=12$  CA mice). **e<sub>i</sub>**, After lick-training, mice underwent auditory cued go/no-go training that consisted of ~300 trials per session. **e<sub>ii</sub>**, Mice learn to lick after a S+ tone to obtain a plain water reward (3ul) and withhold licking to an S- tone to avoid a time-out. **e<sub>iii</sub>**) Correct responses are hits (licking to the S+ tone) and correct rejects (withhold licking to the S-), while incorrect responses are false alarms (licking to the S-) and misses (not licking to the S+). **f**, Accuracy comparison between WR85 (blue) and CA (yellow) mice on highly engaged trial blocks is similar (group comparison ANOVA,  $F(1,20)=3.06$ ,  $p=0.081$ , interaction group x trials  $F(1,20)=0.83$ ,  $p=0.68$ ) ( $n=11$  WR85, and  $n=12$  CA mice). Expert accuracy is defined as 75% correct (gray horizontal line). **g**, No differences were observed in the number of trials to reach expert accuracy between groups (Wilcoxon rank sum test,  $p=0.42$ ) ( $n=11$  WR85, and  $n=12$  CA mice).

129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152

153 **Abrupt transitions from goal-directed to habitual behavior spontaneously occur**  
154 **during training**

155 Mice were then trained on a discriminative auditory go/no-go task in which they learned  
156 to lick to one tone (S+, reward-predicting cue) for a water reward (hit) and withhold licking  
157 to another tone (S-, non-reward-predicting cue) to avoid a timeout (**Fig. 1e**, correct reject).  
158 CA mice exhibited lower response rates to the S+, consistent with the reduced motivation,  
159 but surprisingly only minor differences in discrimination performance throughout learning.  
160 Specifically, when restricting performance assessment to blocks of high engagement  
161 (>50 trials with >90% hit rate), task performance during learning was similar to WR85  
162 mice (**Fig. 1f-g**). In addition, CA and WR85 mice exhibited similar, high discrimination  
163 performance (75%) within 1,500 trials (WR85=1132 ± 87 and CA=1422 ± 234 trials,  
164 Wilcoxon rank-sum test p=0.93). Overall, these data show that CA mice learn task  
165 contingencies at similar rates to WR85 mice while responding less to reward-predicting  
166 cues.

167 We next sought to explore in detail the impact of reduced motivation on responses  
168 to the reward-predicting cue. This could be driven by a continuously lower response rate  
169 (generally lower motivation) or, alternatively, a fluctuating response rate (periodic  
170 changes in motivation). In WR85 mice, animals initially increase their action rates to both  
171 tones, followed by a slow reduction in response to the S- (**Fig. 2a**, left example animal).  
172 The S+ response (hit rate) stayed consistently high with minimal variability. We observed  
173 a striking contrast in CA mice (**Fig. 2a**, right), where for thousands of trials, CA mice  
174 exhibit a fluctuating hit rate, regularly shifting from epochs of high hit rate to low hit rate,  
175 suggesting that mice are volitionally engaging in the task and that we can track their  
176 fluctuating motivation levels. We thus focused on the hit rate as the response rate of  
177 interest (**Fig. 2b**). CA mice showed significantly fewer blocks of high engagement (hit rate  
178 > 90%) compared to WR85 mice (**Fig. 2c**, Wilcoxon rank sum test, p=0.0028).  
179 Surprisingly, after this high-variability phase, most CA mice abruptly transitioned to a low-  
180 variability phase (**Fig. 2d**, red line). We observed that 10 out of 12 (83.3%) of the CA  
181 mice exhibited high hit rate variability, of which 8 out of 10 (80%) transitioned to a low-  
182 variability phase (**Fig. 2e**). Transitions occurred more than 2,000 trials after animals  
183 exhibited expert discrimination performance, suggesting this transition was not due to



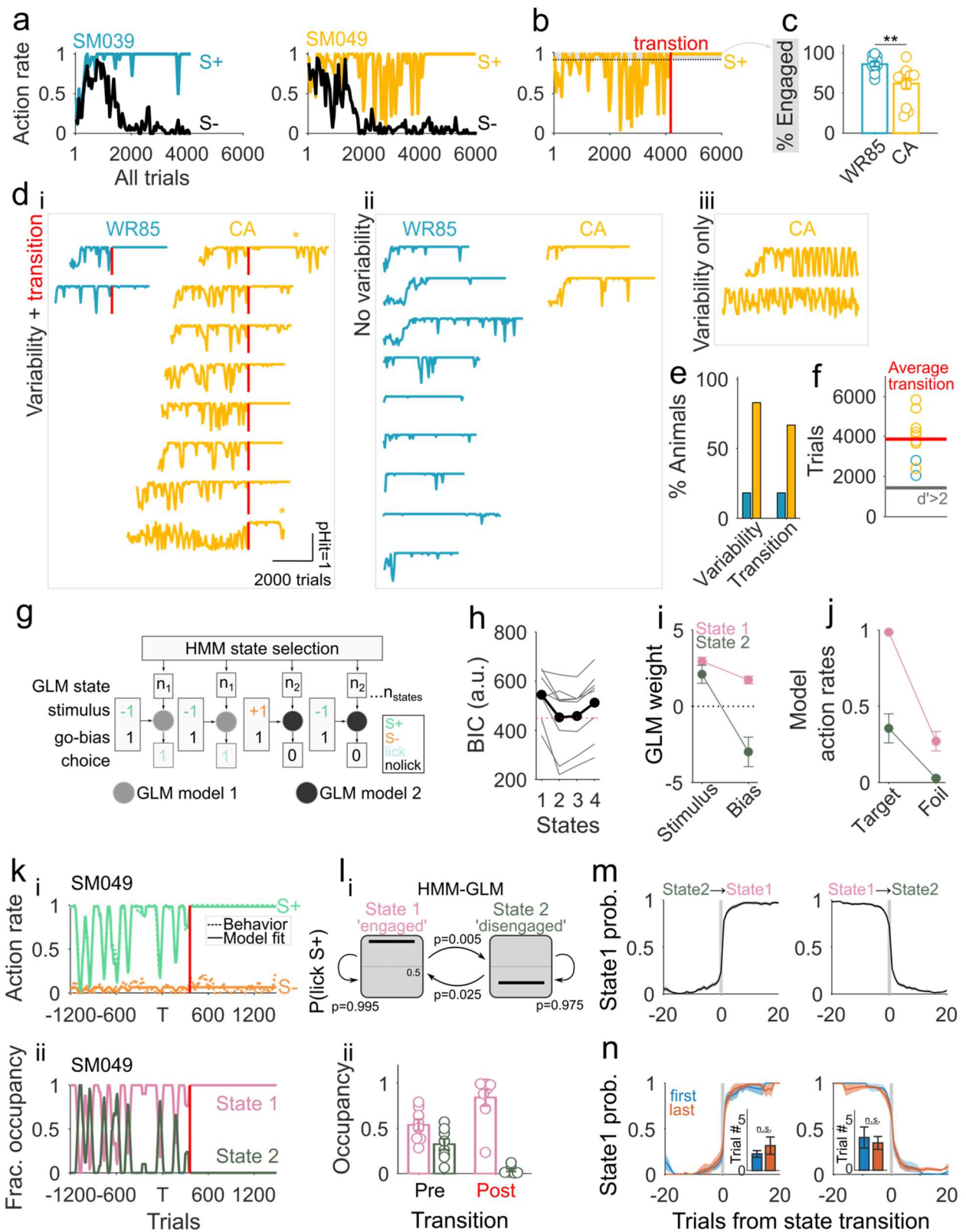
184 ongoing contingency learning ( $d' > 2$  expert=1436±454 trials vs. transition=3832±385  
185 trials) (**Fig. 2f**). To confirm that this change in hit rate was not due to a sudden change in  
186 underlying motivation, we measured the weights of the animals daily. Importantly, we  
187 observed (1) no differences in discrimination performance around the transition (paired t-  
188 test,  $p=0.83$ ) (**Extended Data Fig. 2a**) (2) no evidence of changes in weight pre- versus  
189 post-transition (**Extended Data Fig. 2b**) (paired t-test,  $p=0.19$ ), and (3) no changes in  
190 consumption in the home cage based on measurements of post-task and next-day  
191 weights (**Extended Data Fig. 2c**) (paired t-test,  $p=0.81$ ). This suggests that CA mice do  
192 not exhibit increased motivation post-transition and, instead, points to the rapid  
193 emergence of habitual control. Interestingly, we observed this transition typically occurred  
194 at the beginning of a new session (**Extended Data Fig. 2d-f**), suggesting that transitions  
195 from goal-directed to habitual behavior may be supported by offline processing.

196 Our analysis thus far required categorization of behavior based on pre-defined  
197 criteria (low versus high variability, pre- versus post-transition) and experimenter-defined  
198 parameters (**Extended Data Fig. 3a-d**). We sought to test whether a bottom-up, model-  
199 based approach could identify behavioral 'states' in an unbiased, and trial-by-trial,  
200 manner. To do this, we applied a generalized linear model that incorporates a hidden  
201 Markov process (HMM-GLM)<sup>6</sup> on trial-by-trial behavioral data after animals reached  
202 expert discrimination performance (**Fig. 2g**). The HMM-GLM identified two states that  
203 best described the behavior in expert animals, as defined by the lowest cross-validated  
204 Bayesian Information Criterion value (BIC) (**Fig. 2h-i**). Both states were sensitive to the  
205 stimulus but with distinct action biases. State 1 (pink) exhibited high engagement, evident  
206 by a high bias and high hit rate (i.e. action rate on target trials), while State 2 (dark green)  
207 exhibited a strong disengagement, evident by a low bias and low hit rate (**Fig. 2i-j**). The  
208 HMM-GLM (solid line) accurately recapitulated the behavioral data (**Fig. 2k<sub>i</sub>**, and  
209 **Extended Data Fig. 3e**). Interestingly, before the transition, CA mice regularly switched  
210 between the two states both within and across sessions. After the habitual transition,  
211 however, State 1 ('engaged') dominated behavioral performance (**Fig. 2k<sub>ii-l</sub>**, and  
212 **Extended Data Fig. 3e**). We then used the HMM-GLM model to predict the transition in  
213 a bottom-up manner which we found to be similar to the behaviorally predicted one while  
214 providing greater temporal specificity (**Extended Data Fig. 3f**) (Wilcoxon rank sum test,

215 p=0.74). The model-defined shifts from ‘Engaged’ to ‘Disengaged’ states before the  
216 habitual transition were strikingly abrupt (**Fig. 2m**). Moreover, the last transition, reflecting  
217 the putative transition between goal-directed and habitual behavior, was as abrupt as the  
218 earlier fluctuations (**Fig. 2n**, ‘last’). These state transitions occurred within less than 5  
219 trials, and for most animals, they occurred at the very beginning of a session (**Extended**  
220 **Data Fig. 3g-i**) further implicating a role for offline processing in the transition of  
221 behavioral control. Thus, both quantification of behavioral data and model-based  
222 approaches using the HMM-GLM converge on the abruptness of the identified habitual  
223 transition (**Extended Data Fig. 3a-f**).

224





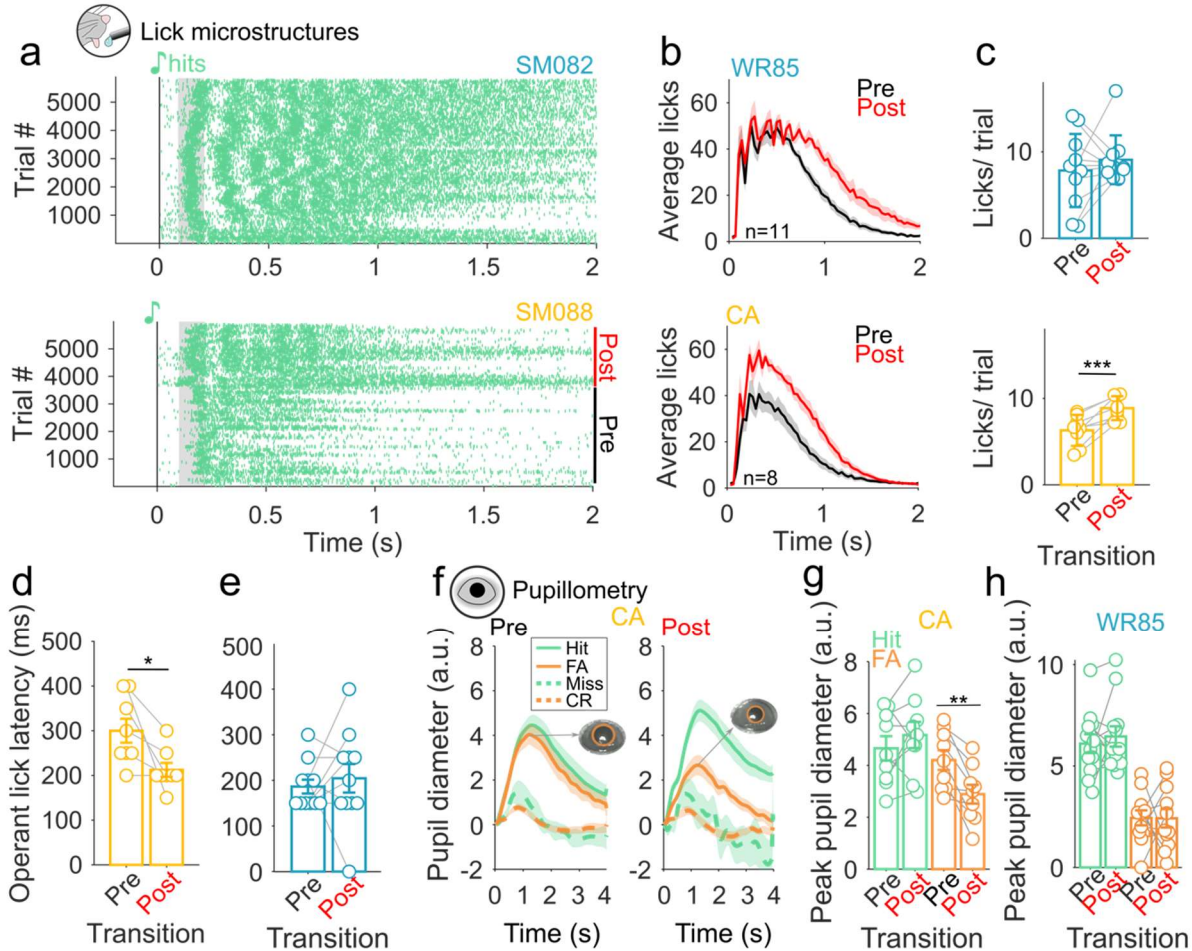
226 **Figure 2. Abrupt state-like transitions from goal-directed to habitual behavior appear spontaneously**  
227 **within individual animals. a**, Hit and FA action rates for example WR85 (left) and a CA (right) animals. **b**,  
228 Hit rate of CA example animal (in a) showing periods of high (gray shadow) and low engagement, followed  
229 by a spontaneous transition (red vertical line) to low hit rate variability. **c**, CA mice have significantly less  
230 engaged blocks compared to WR85 (Wilcoxon rank-sum test  $p=0.0028$ ). **d<sub>i</sub>**, Most CA mice (8/12, 66.7%)  
231 showed hit-rate variability and the presence of a transition, while only a low percentage of WR85 mice did  
232 (2/11, 18.2%). Some CA mice that transitioned to low variability hit rate, seemed to transition to high  
233 variability after a while (asterisks). **d<sub>ii</sub>**, Most WR85 mice (9/11, 81.8%) showed no hit rate variability, while  
234 only a few CA mice did (2/12, 16.7%). **d<sub>iii</sub>**, Few CA mice showed initial hit rate variability, but never  
235 transitioned to low variability (2/12, 16.7%). **e**, Overall, most CA mice (10/12, 83.3%) showed high action  
236 rate variability, while only 18.2% of WR85 did. A high percentage of CA mice also showed transitions  
237 (66.7%), while only 18.2% of WR85 mice did ( $n=11$  WR85 mice and  $n=12$  CA mice). **f**, From all animals  
238 that transitioned, the average transition session is Session 13, which occurs thousands of trials after  
239 reaching expert performance (session 6) ( $n=11$  WR85 mice and  $n=12$  CA mice). **g**, An HMM-GLM was  
240 used to model behavioral data, which allows us to analyze the state-like nature of transitions. **h**, An HMM-  
241 GLM with two states provides the best fit for most of the CA animals based on a BIC analysis ( $n=8$  CA  
242 mice). **i**, Both states are equally stimulus-driven, but state 2 is characterized by a NoGo, or disengaged  
243 bias ( $n=8$  CA mice). **j**, State 1 (pink) is highly stimulus selective between target and foil trials with high  
244 engagement, while State 2 (green) shows overall task disengagement ( $n=8$  CA mice). **k**, A CA exemplar  
245 shows that the two-state HMM-GLM model accurately recapitulates the behavior (top), and the two states  
246 govern the pre-transition phase, while only one state becomes explanatory of the post-transition phase. **l**,  
247 The GLM states reflect transitions between an engaged state (state 1, pink, hit rate = 0.98) and a  
248 disengaged state (state2, olive, hit rate = 0.35). Across all mice, the HMM-GLM model predicted that the  
249 probability of staying in engaged or disengaged state (trial-by-trial) is 0.995 and 0.975 respectively, whereas  
250 the transition probability between states is 0.005 (engaged to disengaged) and 0.025 (disengaged to  
251 engaged). **l<sub>ii</sub>**, State occupancy pre-transition (black) is approximately divided 50%-50% between State 1  
252 and State 2, while post-transition (red), State 1 dominates ( $n=8$  CA mice). **m**, NoGo to Go (State 2 → State  
253 1) transitions and Go to NoGo (State 1 → State 2) transitions happening in the goal-directed phase, are  
254 abrupt ( $n=8$  CA mice). **n**, We observed no differences between the first (blue, belonging to the goal-directed  
255 phase), and last (red, belonging to goal-directed to habitual) transitions in abruptness. Both happen within  
256 1-4 trials ( $n=8$  CA mice).

257  
258

## 259 **Licking microstructures demonstrate automaticity post-transition**

260 One alternative interpretation of these results is that animals abruptly transitioned to a  
261 higher vigilance state in which they exploit their task knowledge in a goal-directed  
262 manner, rather than a transition to habitual decision-making. In skill-based learning, motor  
263 patterns of ‘automaticity’ can be used as evidence for the formation of habits<sup>29</sup>. We  
264 analyzed lick microstructures in detail to determine the extent to which transitions in hit-  
265 rate variability were concomitant with changes in motor automaticity. Before transitions,  
266 CA mice exhibited highly variable lick microstructures in comparison to WR85 mice (**Fig.**  
267 **3a**, top). Post-transition, however, three aspects of their licking behavior abruptly appear:  
268 a uniform lick stereotypy (**Fig. 3a**, bottom), an increase in consummatory licks (**Fig. 3b-**  
269 **c**, bottom), and a reduction in reaction time (**Fig. 3d**). These patterns were consistent  
270 across trials and sessions after the habitual transition demonstrating the simultaneous  
271 appearance of motor automaticity.

272



273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

**Figure 3. Motor automaticity and error-related pupillary signatures appear concomitant with habit transitions.** **a**, Exemplar lick raster plots of a WR85 (upper) and a CA (lower) animal, showing individual licks (green lines) to target tones throughout training. Tone onset (0, black note) is followed by a dead period (100ms) and the presence of operant licks (gray rectangle). **b**, CA mice show a strong increase in post-transition number of consummatory licks (bottom) compared to WR85 mice (top). **c**, The average number of licks per trial is significantly higher in CA mice post-transition compared to pre (bottom,  $p=0.0019$ ) in comparison with WR85 mice (top, paired t-test,  $p=0.41$ ) ( $n=11$  WR85 mice and  $n=8$  CA mice). **d**, A significant reduction in the operant lick latency is observed post-transition (paired t-test,  $p=0.016$ ) ( $n=8$  CA mice). **e**, No changes in operant lick latency for WR85 mice (paired t-test,  $p=0.59$ ) ( $n=11$  WR85 mice). **f**, Evoked pupil dilation is significantly reduced for false alarms (orange) post-transition compared to evoked responses during hits (green) (paired t-test,  $p=0.0019$ ) ( $n=9$  recording days in total of  $n=3$  CA mice for pre and post-transition). **g**, A significant difference was observed in FA evoked pupil dilation between pre and post transition in CA mice (non-parametric paired t-test,  $p=0.0039$ ) but not for hits (non-parametric paired t-test,  $p=0.25$ ) ( $n=9$  CA datapoints, from 3 mice, corresponding to 3 days pre and 3 days post transition). **h**, No differences were observed in tone evoked pupil dilation during FA (non-parametric paired t-test,  $p=0.90$ ), or Hits (non-parametric paired t-test,  $p=0.99$ ) between pre and post transition in WR85 mice ( $n=12$  WR85 4 mice, corresponding to 3 days pre and 3 days post transition).

## 295 **The underlying decision process is reflected in pupillary dynamics**

296 With the ability to pinpoint the transition from goal-directed to habitual behavior, we next  
297 sought to examine whether the decision controller being used could be inferred from  
298 changes in pupillary dynamics. When behavior is under goal-directed control, the  
299 execution of an action is driven by the expectation of reward (the action-outcome  
300 contingency). In contrast, when behavior is under habitual control, the execution of the  
301 action is driven by the presence of a cue (the stimulus-action contingency) and errors are  
302 likely due to ‘slips of action’ with little relation to the outcome expectation<sup>30</sup>. One tool  
303 commonly used as a biomarker of decision-making processes is the pupillary response,  
304 as pupil dynamics reflect choice<sup>31</sup> and track value-based decision-making<sup>32</sup>. Changes in  
305 pupil size are associated with increased reward magnitude, task investment or effort<sup>33,34</sup>,  
306 and reward expectation<sup>35</sup>. Here, we hypothesized that habitual behavior should recruit  
307 less cognitive effort and lower sensitivity to errors and would be reflected in changes in  
308 pupillary dynamics.

309 To test this, we recorded and quantified trial-level pupil responses in a subset of  
310 CA (n=3) and WR85 (n=4) mice. We aligned our phasic, task-evoked pupil measurements  
311 to the transition trial (empirically calculated and confirmed with the HMM-GLM) in CA  
312 mice. We observed that consistent with previous work<sup>31,36</sup>, false alarms elicited strong  
313 pupil dilations (expert and pre-transition, **Fig. 3f**, left and **Fig. 3g**, left). During habitual  
314 behavior (post-transition), however, false alarms elicited a much weaker pupil dilation  
315 (**Fig. 3f**, right and **Fig. 3g**, right). These changes could not be explained by commonly  
316 reported movement-evoked pupillary response<sup>37–40</sup> as the rate of false alarms (**Extended**  
317 **Data Fig. 4a**) and licks per false alarm (**Extended Data Fig. 4b**) were similar pre- and  
318 post-transition. This effect was not observed in WR85 mice as the underlying motivational  
319 state likely overwhelms subtle changes in pupillary dynamics (**Fig. 3h**, and **Extended**  
320 **Data Fig. 4c**). These data demonstrate that the differences in tone-evoked pupil dilation  
321 reflect decreased error sensitivity during habitual behavior, suggesting that behavioral  
322 control has shifted to a less deliberative and less cognitively demanding controller.

323

324

325



## 326 **Reward devaluation confirms timing of habit transitions**

327 A current gold standard for assessing whether a behavior is under goal-directed or  
328 habitual control is the use of discrete, outcome devaluation sessions<sup>28</sup>. To test whether  
329 our volitional engagement paradigm is consistent with this method, we reasoned that all  
330 animals are likely to transition to habitual behavior but that the transition in WR85 mice is  
331 masked by their continuously high hit rates due to ceiling levels of motivation. We inferred  
332 the transition in WR85 mice using our median transition session from CA mice (session  
333 13). We used discrete satiety test sessions before and after this inferred transition. In  
334 these sessions, mice had *ad libitum* access to plain water for 10 minutes (**Fig. 4a<sub>i</sub>**) prior  
335 to performing the go/no-go task (60 trials). Under goal-directed control, behavior is  
336 expected to be sensitive to satiety (**Fig. 4a<sub>ii</sub>**, left) while under habitual control, behavior is  
337 thought to be insensitive to satiety (**Fig. 4a<sub>ii</sub>**, right). We found that pre-transition (10<sup>th</sup>  
338 session), WR85 mice were highly sensitive to reward devaluation, abolishing responses  
339 (**Fig. 4b**) (Session 10, non-devalued vs devalued paired t-test,  $p=0.0000019$ ), while post-  
340 transition, these same animals were less sensitive (**Fig. 4b**, devalued session 10 vs  
341 devalued session 15, paired t-test  $p=0.0072$ ). Importantly, we observed no differences in  
342 non-devalued (ND) actions rates between pre- and post-transition (**Fig. 4b**, ND session  
343 10 vs ND session 15 paired t-test  $p=0.25$ ). This effect was independent of water  
344 consumption during the devaluation sessions (**Extended Data Fig. 5a**). To further  
345 validate these results obtained in WR85 mice, we tested within-session satiety in CA mice  
346 (satiety-based devaluation was not possible in CA mice, as they only intermittently drank  
347 plain water even when it was freely available, similar to lick training, **Fig. 1**). During goal-  
348 directed behavior (session 10), CA mice exhibited reductions in hit rate when comparing  
349 the beginning to the end of the session (**Fig. 4c**, black) (paired t-test  $p=0.038$ ), suggesting  
350 a session-level impact of satiety. Interestingly, during habitual behavior (session 15), we  
351 observed no such reduction in hit rate (**Fig. 4c**, red) (paired t-test,  $p=0.24$ ), suggesting no  
352 impact of session-level satiety. These data help to validate the volitional engagement  
353 paradigm as a means to assess the underlying decision process.

354

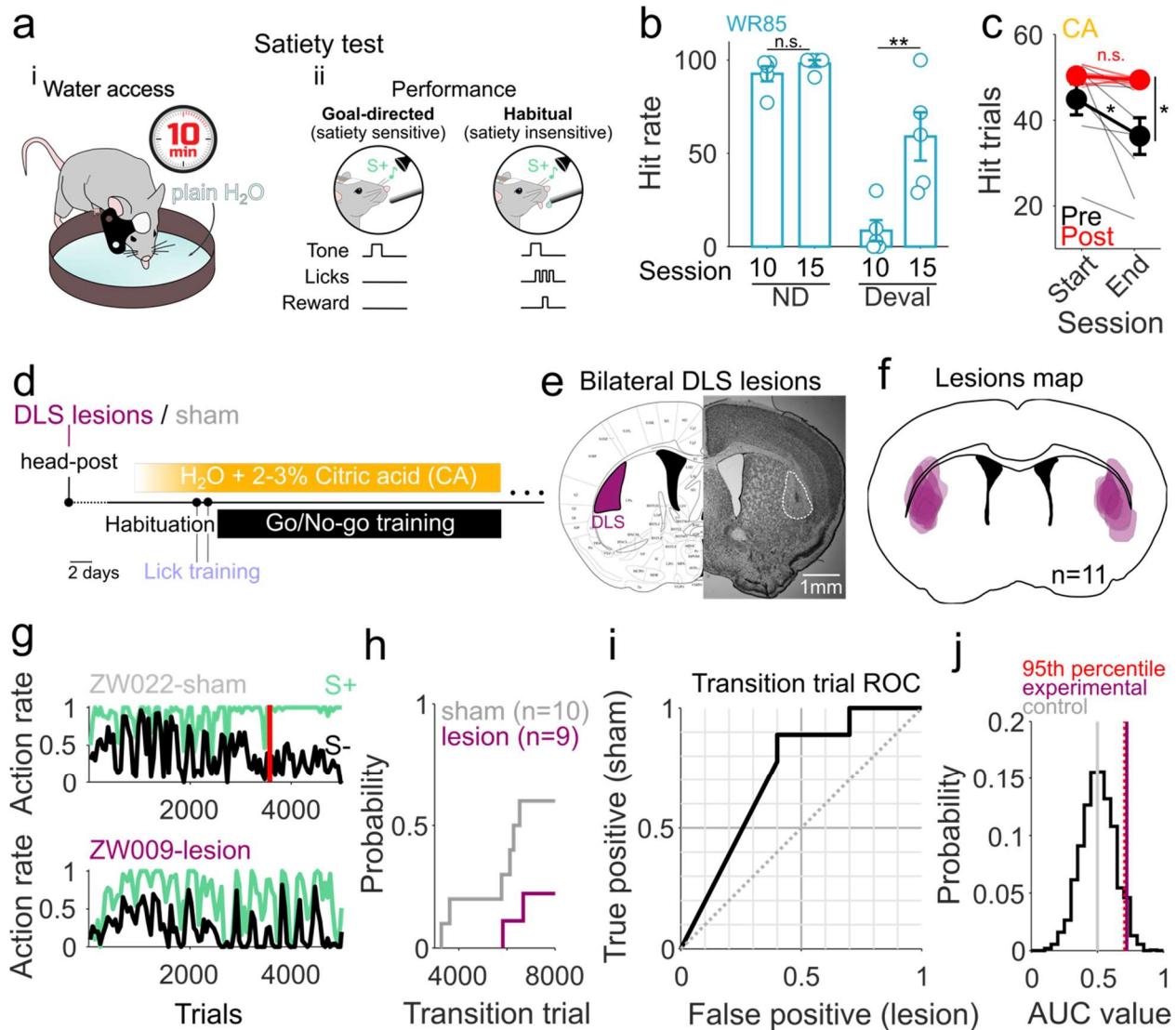
355

356

357 **Bilateral lesions to the DLS prevent transitions to habitual behavior**

358 The two decision processes (goal-directed and habitual) are thought to be sub-served by  
359 distinct neural circuits in the dorsal striatum<sup>41</sup>. The dorsomedial striatum (DMS) and  
360 dorsolateral striatum (DLS) enable goal-directed and habitual behavior, respectively<sup>42,43</sup>.  
361 Using standard outcome devaluation procedures, rodents with lesions to the DLS persist  
362 in goal-directed mode (i.e. when they should be sensitive to outcome devaluation) even  
363 after significant amounts of overtraining<sup>7</sup>. Thus, DLS lesions provide a powerful and  
364 orthogonal approach to test the validity of our volitional engagement paradigm in  
365 assessing the underlying decision mode. To do this, we bilaterally lesioned the DLS in  
366 CA mice (NMDA 20µg/µl, 100nl/site) before head post implantation and behavioral  
367 training (**Fig 4c**). All DLS-lesioned CA mice had visible, localized, and overlapping lesions  
368 (**Fig. 4d-e**, and **Extended Data Fig. 5b**), while shams did not (**Extended Data Fig. 5c-**  
369 **d**). DLS-lesioned CA mice exhibited high variability in hit rates that persisted for much  
370 longer than in sham CA mice (**Fig 4f**, and **Extended Data Fig. 5e**) and lesioned animals  
371 rarely transitioned to habitual behavior (**Fig. 4g-i**, and **Extended Data Fig. 5e**) (60%  
372 sham vs 20% lesioned). Importantly, CA mice with DLS lesions exhibited no significant  
373 deficits in the learning of the task contingencies (**Extended Data Fig. 5f-g**). These data  
374 offer independent evidence, through the manipulation of habit-relevant striatal circuits,  
375 that the transitions we observe are indeed genuine transitions from goal-directed to  
376 habitual behavior.





377

378

379 **Figure 4. Bilateral DLS lesions delay or block transition to habitual behavior.** a, Satiety test. WR85  
 380 mice receive free access to plain water for 10 minutes right before training. a<sub>ii</sub>, During training, if goal-  
 381 directed, mice are expected to reduce their action rates since they are sensitive to satiety. In habitual mode,  
 382 mice tend to maintain a high action rate since they have developed a habit and become less or insensitive  
 383 to satiety. b, Hit rates during Non-devalued (ND, sessions 9) and devalued (Deval, session 10) trials for  
 384 WR85 mice. Pre-transition (session 10) mice show high sensitivity to devaluation (ND vs Deval Session 10,  
 385 paired t-test,  $p=0.0000019$ ). Hit rates during Non-devalued (ND, session 14) and devalued (Deval, session  
 386 15) trials for WR85 mice. Post-transition (session 15) mice show reduced sensitivity to devaluation (ND vs  
 387 Deval Session 15, paired t-test,  $p=0.017$ ), with action rates significantly different from the pre-transition  
 388 phase (paired t-test,  $**p=0.0072$ ) ( $n=5$  WR85 mice). c, Intra-session satiety is observed in CA mice pre-  
 389 transition but not post-transition. While animals pre transition significantly reduce their licking to S+ by the  
 390 end of a session (indicating they have been satiated), post transition, these same animals show no  
 391 differences. The number of target trials with licks is significantly different between the start and end of a  
 392 session pre-transition (black, paired t-test,  $p=0.038$ ), but not post-transition (red, paired t-test,  $p=0.24$ ). A  
 393 significant difference is also observed between pre- and post-transition hit trials at the end of a session  
 394 (end, paired t-test,  $p=0.017$ ), while there are no differences at the start of a session (start, paired t-test,  
 395  $p=0.15$ ) ( $n=8$  CA mice). d, Lesion protocol. Mice undergo an NMDA lesion or sham right before head-post  
 396 implantation. The rest of the water restriction, habituation and training protocol is as shown in Figure 1a. e,

397 Lesion exemplar. Coronal view of the bilateral lesion sites at the DLS. **f**, Lesion map of all animals (n=11  
398 mice), showing homogeneous and overlapping lesions. **g**, Action rate exemplars for sham (top, gray) and  
399 lesioned (bottom, purple) animals. **h**, Cumulative distribution of transition trial for animals that showed  
400 behavioral variability (n=10 sham and n=9 lesioned mice). While 60% of sham animals show a transition,  
401 only 20% of lesioned animals do. **i**, ROC curve built from transitions probabilities in **g**. **j**, AUC values for  
402 shuffled labels in our experimental groups. The difference in AUC between sham and lesioned mice (purple  
403 line) falls into the 95<sup>th</sup> significance percentile (red dotted line), compared to the difference between control  
404 animals (gray line).

405

## 406 **Discussion**

407 A fundamental tenet of animal behavior is the existence of multiple ‘controllers’ that  
408 govern decision-making. One prevailing framework is that instrumental decisions come  
409 about from two distinct processes: goal-directed and habitual<sup>44</sup>. In rodent-based learning  
410 paradigms, goal-directed behaviors are thought to become habitual upon overtraining<sup>45–</sup>  
411 <sup>47</sup>. The goal-directed system dominates early in learning when animals will ‘work’ for food  
412 or water because the outcome itself is desirable. This requires both a representation of  
413 the action-outcome contingency and the recognition of the outcome as a motivational  
414 incentive. When under goal-directed control, behavioral decisions are ‘value-based’ and  
415 flexible but also cognitively demanding<sup>48</sup>. The habitual system is thought to take over  
416 during overtraining to simplify the decision process and reduce cognitive complexity by  
417 shifting to a stimulus-response mode of behavior<sup>48</sup>. When under habitual control,  
418 behavioral decisions can be considered ‘value-less’ and inflexible but also less cognitively  
419 demanding<sup>49</sup>. Over the past 50 years, behavioral, neural, and theoretical support for these  
420 two distinct decision processes (but also the complexity of their interaction) has grown  
421 largely due to behavioral manipulations, including outcome devaluation, contingency  
422 degradation, and the use of different reinforcement schedules. These behavioral tools  
423 have been invaluable to gain a deeper understanding of the behavioral, neural, and  
424 theoretical basis of the multiple systems controlling decision-making. Nevertheless, the  
425 extent to which discrete measures of sensitivity to outcome devaluation sufficiently  
426 distinguishes goal-directed from habitual control is still under scrutiny<sup>50,51</sup> as sensitivity to  
427 outcome devaluation can also be triggered by unexpected cues<sup>52</sup> and in situations where  
428 habits are expected to form<sup>53</sup>. More broadly, the current methodologies remain  
429 fundamentally limited in their temporal resolution and individual specificity, limiting the  
430 assessment of nature, timing, and properties of habit formation<sup>19,50</sup> in individual animals.

431 As a result, an essential question first posed by Clark Hull in 1943 has remained  
432 unresolved: ‘Is this transition abrupt, or is it gradual and progressive?’

433 Here, we lay out an approach that allows real-time assessment of the underlying  
434 decision process without the need for discrete testing sessions and without the  
435 implementation of specified training schedules to bias decision modes<sup>26,54,55</sup>. Rather than  
436 binarize motivation into motivated (non-devalued) and un-motivated (devalued), we  
437 sought to mimic naturalistic motivation levels. We hypothesized that by shifting an  
438 animals’ desire for an outcome from a *need* (survival) to a *preference* (palatability),  
439 animals would gain agency on their motivation to engage in a task based on the desire  
440 for an outcome (in our case, plain water droplets). Here, we show one approach to such  
441 a ‘volitional engagement’ paradigm by giving animals *ad libitum* access to CA water<sup>4,5</sup> in  
442 the home cage, which reduces the motivation for plain water (**Fig. 1b-d**). In goal-directed  
443 mode, animals volitionally engage and disengage from the task, reflected in the behavior  
444 as state-like switching early in learning (see **Fig. 2a** and **Fig. 2d**). As behavior becomes  
445 habitual, animals shift to constant engagement, behaviorally observed as an abrupt shift  
446 to a constantly high action rate (**Fig. 2a**, **Fig. 2d**, **Fig. 2m**, and **Extended Data Fig. 3a-**  
447 **e**), contradicting the assumption that habit expression is gradual. These transitions  
448 occurred thousands of trials after reaching expert discrimination performance (**Fig. 2f**) but  
449 at different time points for individual animals.

450 We then used orthogonal measurements of motor automaticity, pupillary  
451 dynamics, sensitivity to traditional outcome devaluation, and lesions of the DLS to confirm  
452 that our observations reflected a transition to habitual decision-making versus differences  
453 in vigilance or discrimination ability. In other words, behavioral automaticity and reduced  
454 error sensitivity occurs concomitant with habit formation in our paradigm. While  
455 automaticity alone might be a reductionist perspective in habit formation<sup>53</sup>, the composite  
456 picture across behavioral and neural approaches in our study points toward the habitual  
457 nature of behavior post-transition in volitionally engaged mice. This novel approach—  
458 which we term ‘volitional engagement’—adds a powerful *en passant* tool to study habit  
459 formation and perseverance.

460 Pinpointing the precise nature of this transition under naturalistic motivational  
461 conditions provides a powerful tool to explore the psychological and neural basis of habit

462 formation in real-time. An abrupt ‘insight-like’ transition suggests that a higher-level  
463 process operates in conjunction with lower-level associative processes. Our findings  
464 challenge the notion that habit expression is cumulative, with its likelihood increasing  
465 incrementally with each successive reinforcement. We demonstrate that habits appear  
466 nearly instantaneously (within 5 trials) with animals having received vastly different levels  
467 of reinforcement during training. This suggests that a separate higher-level process  
468 arbitrates between goal-directed and habitual control. Factors such as cognitive demand  
469 or environmental uncertainty likely contribute to when the commitment emerges to solve  
470 the task in a simple and inflexible manner, activating an otherwise dormant habitual  
471 controller. The habit becomes instantiated precisely when animals choose to use the  
472 habitual controller, whether they explicitly want the water or not. In this view, animals  
473 under habitual control can still internally ‘experience’ periods of goal-directedness (i.e.,  
474 they still want plain water sometimes). In addition, the higher-level nature of the choice  
475 suggests that habitual control need not be permanent. Interestingly, some animals  
476 reverted to goal-directed mode after several sessions in habit mode (see **Fig. 2d**, yellow  
477 asterisks and **Extended Data Fig. 5e**, asterisks), suggesting that transitions to habitual  
478 decision-making are not intrinsically permanent. This higher-level process that controls  
479 the switch may also help explain why techniques such as outcome devaluation yield  
480 conflicting results in rodents<sup>56</sup> (due to individual variability in when transitions occur) and  
481 remain largely ineffective in humans<sup>19</sup> (given a more complex interaction between  
482 cognitive and motivational drivers).

483 The current consensus, though contentious<sup>57</sup>, is that habitual and goal-directed  
484 behavior are supported by the DLS and DMS, respectively<sup>42</sup>, but studies utilizing discrete  
485 satiety test sessions or experimenter-defined overtraining periods yield confounding and  
486 even conflicting results: some evidence argues that DLS activity changes rapidly before  
487 the behavior onset of a habit<sup>58</sup>, with others finding that the change is more gradual and  
488 closely aligned with a behavioral change<sup>59</sup>. Recent reports even observe an eroding  
489 distinction in the control of actions between the DLS and DMS as training progressed<sup>57</sup>.  
490 The spontaneous and abrupt appearance of habitual control provides a behavioral marker  
491 upon which to identify ‘switch-like’ activity in the underlying neural circuits<sup>19</sup>. The  
492 possibility of a higher-level process that arbitrates between goal-directed and habitual

493 control points to regions such as premotor or prefrontal cortical areas<sup>17</sup>, which have  
494 bidirectional interactions with the DMS and DLS. Alternatively, this arbitration process  
495 may be governed in the striatum itself. Identifying the neural circuit dynamics that govern  
496 this transition remains an important area for future investigation.

497         Finally, our discovery of abrupt transitions to habitual behavior may inform distinct  
498 interventional strategies in habit disorders in humans. Rather than relying on gradual  
499 exposure and/or systematic desensitization, there might be value in interventions that  
500 combine such associative strategies with cognitive control strategies. Our data suggest  
501 that it may be possible to predict when transitions will occur and if such predictions are  
502 possible and can be extended from rodents to humans, it could provide a powerful tool to  
503 interfere or manipulate the emergence of maladaptive habits.

## 504 **Methods**

### 505 **Animals**

506 All mice were housed in standard plastic cages with 1-4 littermates and kept in a 12-h/12-  
507 h light/dark cycle (10:30 am / 10:30 pm) with controlled temperature (19.5-22°C) and  
508 humidity (35-38%). All the mice used in this study were male C57BL/6J from Jackson  
509 Labs (strain 664) with an age of  $11.61 \pm 0.21$  (average  $\pm$  SEM) weeks at the start of  
510 training). All the experimental and surgical procedures were approved and performed in  
511 accordance with the Johns Hopkins University IACUC protocol (license # MO20A272).

512

### 513 **Surgical procedures**

514 Mice were anesthetized with isoflurane (5.0% at induction, 1.5-2.5% during surgery) and  
515 placed on a stereotactic apparatus (Kopf). The hair over their skull was removed with hair  
516 removal cream and the area disinfected with betadine. The skin over the skull was  
517 removed and the area was cleaned of connective tissue with 3% H<sub>2</sub>O<sub>2</sub>. A custom-made  
518 stainless-steel head-post was fixed onto the exposed skull with C&B Metabond dental  
519 cement (Parkell). The animals were given 1-3 days to recover. Mice that underwent  
520 bilateral DLS excitotoxic lesions or sham injections, received NMDA (Sigma Aldrich,  
521 20 $\mu$ g/ $\mu$ l NMDA in PBS1x with 10% glycerol) or vehicle (PBS1x with 10% glycerol)  
522 respectively (with a Hamilton syringe and a Harvard Apparatus Pump 11 Elite,  
523 100nL/injection-site at a 70nL/min). The injections were made at AP+1.0mm, DV $\pm$ 2.6mm,  
524 ML-2.8mm via burr holes which were sealed with Jet Denture Repair Acrylic (Lang Dental)  
525 prior to headpost implantation.

526

### 527 **Histology**

528 At the end of the DLS lesion experiment, the brains of all animals were obtained via  
529 transcardiac perfusion<sup>60</sup> and stored in 4% paraformaldehyde solution in PBS1x overnight.  
530 After further dehydration in 30% sucrose (Sigma-Aldrich), the brains were frozen in OCT  
531 gel (Tissue-Tek®) and sliced using a cryostat (Leica) into 50 $\mu$ m slices. The slices were  
532 mounted on gelatin-coated slides (FD Neuro) and left in room temperature to dry  
533 overnight. The following day, the slides were stained using 1% cresyl violet (Sigma-  
534 Aldrich) solution and cover glasses were placed and fixated with Permount mounting



535 medium (Fisher Chemical). The slides were imaged under Brightfield settings in a Zeiss  
536 upright microscope (Axio Zoom.V16).

537

### 538 **Habituation and water restriction paradigms**

539 After recovery from surgery, animals were handled and habituated prior to the start of  
540 training for at least 10 days based on previous studies<sup>61</sup>. Head-fixed experimental CA  
541 mice and their littermate controls (WR85) underwent the same surgical, habituation and  
542 testing procedure. Animals were handled by the experimenter/s at increasing times every  
543 day, exposed, and habituated to the head fixation station. The different water restriction  
544 paradigms started after at least 3 days of handling. The standard water restriction (WR85)  
545 protocol prevented the animals from accessing water in their home cage. The mice were  
546 weighed daily, and a limited amount of water (~1.0 mL) was given individually to maintain  
547 80-85% of their original weight. For naturalistic water restriction with citric acid (CA),  
548 animals had *ad libitum* access in their home cage to a bottle of tap water with citric acid  
549 dissolved. The mice were slowly introduced to the taste of CA, increasing its  
550 concentration daily from 0.5% CA to 1-3%, and adjusted accordingly within this range to  
551 keep the animals at ~95% of their original weights.

552

### 553 **Behavioral training**

554 All behavioral training was done using Bpod State Machines (r1 or r2, Sanworks). After  
555 habituation, mice underwent an initial instrumental training phase where they were head-  
556 fixed and trained to lick from a lick tube by rewarding each lick with a drop of water (3  $\mu$ l).  
557 There was no tone stimulus presented during lick training. The lick training session ended  
558 either when 1 ml of water was consumed, or session had reached 30 minutes. On a  
559 subsequent session, mice began training on a go/no-go auditory task. Behavioral events  
560 (trial structure, stimulus and reward delivery, lick detection) were controlled and stored  
561 using a custom-written MATLAB program (2018b, The MathWorks) interfacing with the  
562 Bpod, an electrostatic speaker driver (E1, TDT) and an infrared beam for lick detection.  
563 In a subset of animals, facial movements and pupil size were measured with a Raspberry  
564 Pi (3B) and a Raspberry Pi camera module (NoIR v2) coupled with a Bright-Pi infrared  
565 LED array (PiSupply). Mice were head-fixed inside a Plexiglass tube facing a lick-tube. A

566 free field electrostatic speaker (ES1, TDT) was located ~5 cm from the animal's left ear  
567 and each sound (either 4757 Hz or 8000 Hz, as target or foil stimuli) was calibrated to an  
568 intensity of 60-62 dB (SPL). The pupil camera and IR LED array were positioned ~6 cm  
569 away from the animal's face in a 60-degree angle. Everything was enclosed in a custom-  
570 made sound-attenuated box. Target and foil tones were pseudo-randomly ordered  
571 (equilibrated every 20 trials). Each trial consisted of a pre-stimulus no-lick period (2 s),  
572 stimulus presentation (100 ms), delay (100 ms), response period (2 s) and variable inter-  
573 trial interval (ITI). Typically, mice were trained for ~300-320 trials per with a short block of  
574 20 non-reinforced trials interleaved in the middle of the session<sup>62,63</sup>. Training lasted for a  
575 maximum of 30 days.

576

### 577 **Behavioral analysis**

578 Individual-animal action rates were measured in blocks of 50 trials to obtain hit and false-  
579 alarm rates in discrete but small blocks that allowed us to observe behavioral variability.  
580 Behavioral discriminability was calculated using the z-scored hit rate minus the z-scored  
581 false-alarm rate ( $d'$ ). To avoid infinite values when rates of 1 or 0 are present, the values  
582 were corrected by  $1-1/2N$  or  $1/2N$  respectively, where  $N$  corresponds to the number of  
583 trials. For all the data presented in this paper, we considered animals' to have effectively  
584 learned the task by calculating  $d'$  during non-reinforced trial blocks, which has previously  
585 been demonstrated as an accurate measure of task acquisition<sup>62</sup>. Only mice with a  $d' > 2$   
586 during these non-reinforced trial blocks were included in the analysis (corresponding to  
587 45 out of 48 mice tested, 1 WR85, 1 CA-sham and 1 CA-lesioned mice did not learn the  
588 task and thus were excluded).

589

### 590 **HMM-GLM model implementation**

591 We fit a GLM-HMM model to trial-by-trial choices of each mouse from 4 days before  
592 putative habitual transitions to 4 days after the putative transition (9 days in total). Each  
593 state in HMM contains a Bernoulli GLM defined by a weight vector specifying how  
594 stimulus inputs and bias are integrated in that state. The model was fit using a previously  
595 published expectation-maximization (EM) algorithm<sup>6</sup>. To identify the optimal number of  
596 states, we evaluated the cross-validated BIC by fitting choice data from the 5th day before

597 and after habitual transition. A 2-state model was sufficient to explain the choice behavior  
598 of six animals, capturing an engaged state and a disengaged state, whereas a 3-state  
599 model was needed for two animals, capturing an additional low-discrimination state. For  
600 these two animals, we focused only on the engaged and the disengaged state in  
601 subsequent analysis. To compute state occupancy, we first inferred the behaviorally  
602 dominant state as the state with the highest probability in each trial, and then calculated  
603 the percentage of trials that a state is dominant in a 50-trial bin. We inferred the habitual  
604 transition by identifying the last trial bin where the occupancy of disengaged state was  
605 above a threshold of 30%. The number of trials needed for transitions between engaged  
606 and disengaged state is calculated by the number of trials needed for the dominant state  
607 to reach 75% probability after the transition. To compare the model-inferred transitions  
608 with behavioral data, we quantified the slope of inferred state probability by the GLM (z-  
609 scored) at the trial of state transition, compared to the slope of hit rate changes during  
610 state transitions (z-scored), quantified using various bin sizes around the transition.

611

### 612 **Preprocessing of pupillometry data**

613 20 minutes long pupillometry videos (n=5) were taken as the training dataset for a  
614 DeepLabCut<sup>64,65</sup> (DLC) pre-trained model (resnet\_50). Manual labeling of pupil contour  
615 consisted of 8 points (up, down, left, right, up-left, up-right, down-left, down-right) across  
616 180 randomly selected frames. The network was trained for 564,000 iterations until the  
617 loss rate plateaued. The final network was used to analyze the pupillometry videos from  
618 the experiment. Custom MATLAB code (The MathWorks, 2019b or 2022a) was then used  
619 to remove blink artifact, reconstruct pupil diameter, and apply a low pass filter (3 Hz) to  
620 the data. Individual trials for individual animals were normalized to the median pupil  
621 diameter during correct reject trials per session.

622

### 623 **Statistical Analysis**

624 All analyses were performed using custom-written MATLAB code (The MathWorks,  
625 2019b or 2022a). All datasets were tested for normality using a one-sample Kolmogorov-  
626 Smirnov test; then, parametric or non-parametric statistical tests were applied  
627 accordingly. Two-sample t-tests were used for parametric data, and Wilcoxon rank sum

628 tests were used for non-parametric data. Where required, paired comparisons were  
629 made. For multiple comparison analyses, 2-way ANOVA was performed. To build a  
630 Receiver Operating Characteristic (ROC curve) (**Fig. 4h**) we used the transition  
631 probability of lesioned and sham animals to obtain the area under the curve (AUC) and  
632 generate a shuffled probability distribution to statistically test our experimental animals'  
633 distribution difference. Significance was determined as the difference in AUC value  
634 between lesioned and sham animals when it fell beyond the 95<sup>th</sup> percentile confidence  
635 interval (**Fig. 4i**). All confidence intervals correspond to  $\alpha=0.05$ . Significance is  
636 represented as n.s.  $p>0.05$ , \*  $p\leq 0.05$ , \*\*  $p\leq 0.01$ , \*\*\*  $p\leq 0.001$ , and \*\*\*\*  $p\leq 0.0001$ .

637

### 638 **Data reporting**

639 Sample sizes were determined based on standard cohort sizes from relevant literature.  
640 Mouse allocation to specific groups was randomized but the experimenters were not  
641 blinded to group types.

642

### 643 **Reporting Summary**

644 Further information on research design will be available in the Nature Portfolio Reporting  
645 Summary linked to this article.

646

### 647 **Data availability**

648 Data will be made available upon acceptance of this manuscript.

649

### 650 **Code availability**

651 No specialized software was developed for this work.

652

### 653 **Author information**

#### 654 **Contributions**

655 SM and KVK designed the project. ZW, SM and AL performed the experiments. SM, ZZ,  
656 and RS analyzed the data. ZZ, RS, and AC performed computational modeling. SM  
657 performed final analysis, figures, and data curation. SM, ZW, and KK wrote the

658 manuscript. KVK provided funding and supervised the project. All authors participated in  
659 results interpretation and manuscript editing.

660

661 **Ethics declarations**

662 **Competing interests**

663 The authors declare no competing interests.

664

665 **Acknowledgements**

666 We thank P. Janak, A. Haith, J. Krakauer, N. Kothari, P. Holland, and Y. Cheng for  
667 helpful comments on the manuscript. We thank E. Barker and D. Udzenski for animal  
668 care taking. This work was supported by grants from the NIH R01DC018650 and  
669 R00DC015014 to KVK.

670

671 **Extended data and figures**

672 Extended figures 1 to 5 and legends are provided.

673

674 **REFERENCES**

675

- 676 1. Hull, C. L. Principles of behaviour: an introduction to behavior theory. (Appleton-Century-  
677 Crofts, Inc., 1943).
- 678 2. Daw, N. D. & O'Doherty, J. P. in *Neuroeconomics* 393–410 (Elsevier, 2014).  
679 doi:10.1016/B978-0-12-416008-8.00021-8
- 680 3. Lally, P., van Jaarsveld, C. H. M., Potts, H. W. W. & Wardle, J. How are habits formed:  
681 Modelling habit formation in the real world. *Eur. J. Soc. Psychol.* 40, 998–1009 (2010).
- 682 4. Reinagel, P. Training rats using water rewards without water restriction. *Front. Behav.*  
683 *Neurosci.* 12, 84 (2018).
- 684 5. Urai, A. E. et al. Citric Acid Water as an Alternative to Water Restriction for High-Yield  
685 Mouse Behavior. *eNeuro* 8, (2021).
- 686 6. Ashwood, Z. C. et al. Mice alternate between discrete strategies during perceptual  
687 decision-making. *Nat. Neurosci.* 25, 201–212 (2022).
- 688 7. Yin, H. H., Knowlton, B. J. & Balleine, B. W. Lesions of dorsolateral striatum preserve  
689 outcome expectancy but disrupt habit formation in instrumental learning. *Eur. J. Neurosci.*  
690 19, 181–189 (2004).
- 691 8. Gillan, C. M., Robbins, T. W., Sahakian, B. J., van den Heuvel, O. A. & van Wingen, G. The  
692 role of habit in compulsivity. *Eur. Neuropsychopharmacol.* 26, 828–840 (2016).
- 693 9. Wood, W. & Runger, D. Psychology of Habit. *Annu. Rev. Psychol.* 67, 289–314 (2016).
- 694 10. Yin, H. H. & Knowlton, B. J. The role of the basal ganglia in habit formation. *Nat. Rev.*  
695 *Neurosci.* 7, 464–476 (2006).
- 696 11. Nilsen, P., Roback, K., Brostrom, A. & Ellstrom, P.-E. Creatures of habit: accounting for the  
697 role of habit in implementation research on clinical behaviour change. *Implement. Sci.* 7, 53  
698 (2012).
- 699 12. van Elzelingen, W. et al. Striatal dopamine signals are region specific and temporally stable  
700 across action-sequence habit formation. *Curr. Biol.* 32, 1163–1174.e6 (2022).
- 701 13. Devan, B. D., Hong, N. S. & McDonald, R. J. Parallel associative processing in the dorsal  
702 striatum: segregation of stimulus-response and cognitive control subregions. *Neurobiol.*  
703 *Learn. Mem.* 96, 95–120 (2011).
- 704 14. Smith, K. S., Virkud, A., Deisseroth, K. & Graybiel, A. M. Reversible online control of  
705 habitual behavior by optogenetic perturbation of medial prefrontal cortex. *Proc. Natl. Acad.*  
706 *Sci. USA* 109, 18932–18937 (2012).
- 707 15. Bouton, M. E. *Learning and Behavior: A Contemporary Synthesis.* 556 (Sinauer Associates  
708 is an imprint of Oxford University Press, 2016).
- 709 16. Thorndike, E. L. *Animal intelligence : an experimental study of the associative processes in*  
710 *animals / by Edward L. Thorndike.* (Macmillan,, 1898). doi:10.5962/bhl.title.25848



- 711 17. Lingawi, N. W. & Dezfouli, A. The psychological and physiological mechanisms of habit  
712 formation. *The Wiley handbook on ...* (2016).
- 713 18. Ostlund, S. B. & Balleine, B. W. On habits and addiction: An associative analysis of  
714 compulsive drug seeking. *Drug Discov. Today. Dis. Models* 5, 235–245 (2008).
- 715 19. de Wit, S. et al. Shifting the balance between goals and habits: Five failures in experimental  
716 habit induction. *J. Exp. Psychol. Gen.* 147, 1043–1065 (2018).
- 717 20. Ersche, K. D. et al. Carrots and sticks fail to change behavior in cocaine addiction. *Science*  
718 352, 1468–1471 (2016).
- 719 21. Dickinson, A. Actions and habits: the development of behavioural autonomy. *Philos. Trans.*  
720 *R. Soc. Lond. B, Biol. Sci.* 308, 67–78 (1985).
- 721 22. Adams, C. D. & Dickinson, A. Instrumental responding following reinforcer devaluation. *The*  
722 *Quarterly Journal of Experimental Psychology Section B* 33, 109–121 (1981).
- 723 23. Schreiner, D. C., Renteria, R. & Gremel, C. M. Fractionating the all-or-nothing definition of  
724 goal-directed and habitual decision-making. *J. Neurosci. Res.* 98, 998–1006 (2020).
- 725 24. Vandaele, Y. & Janak, P. H. Defining the place of habit in substance use disorders. *Prog.*  
726 *Neuropsychopharmacol. Biol. Psychiatry* 87, 22–32 (2018).
- 727 25. Holland, P. C. Relations between Pavlovian-instrumental transfer and reinforcer  
728 devaluation. *J Exp Psychol Anim Behav Process* 30, 104–117 (2004).
- 729 26. Gremel, C. M. & Costa, R. M. Orbitofrontal and striatal circuits dynamically encode the shift  
730 between goal-directed and habitual actions. *Nat. Commun.* 4, 2264 (2013).
- 731 27. Balleine, B. W. & Ostlund, S. B. Still at the choice-point: action selection and initiation in  
732 instrumental conditioning. *Ann. N. Y. Acad. Sci.* 1104, 147–171 (2007).
- 733 28. Rossi, M. A. & Yin, H. H. Methods for studying habitual behavior in mice. *Curr Protoc*  
734 *Neurosci Chapter 8, Unit 8.29* (2012).
- 735 29. Aarts, H. & Dijksterhuis, A. Habits as knowledge structures: automaticity in goal-directed  
736 behavior. *J. Pers. Soc. Psychol.* 78, 53–63 (2000).
- 737 30. Du, Y., Krakauer, J. W. & Haith, A. M. The relationship between habits and motor skills in  
738 humans. *Trends Cogn. Sci. (Regul. Ed.)* 26, 371–387 (2022).
- 739 31. Lee, C. R. & Margolis, D. J. Pupil Dynamics Reflect Behavioral Choice and Learning in a  
740 Go/NoGo Tactile Decision-Making Task in Mice. *Front. Behav. Neurosci.* 10, 200 (2016).
- 741 32. Van Slooten, J. C., Jahfari, S., Knapen, T. & Theeuwes, J. How pupil responses track  
742 value-based decision-making during and after reinforcement learning. *PLoS Comput. Biol.*  
743 14, e1006632 (2018).
- 744 33. Bijleveld, E., Custers, R. & Aarts, H. The unconscious eye opener: pupil dilation reveals  
745 strategic recruitment of resources upon presentation of subliminal reward cues. *Psychol.*  
746 *Sci.* 20, 1313–1315 (2009).

- 747 34. van der Wel, P. & van Steenbergen, H. Pupil dilation as an index of effort in cognitive  
748 control tasks: A review. *Psychon. Bull. Rev.* 25, 2005–2015 (2018).
- 749 35. Fröber, K., Pittino, F. & Dreisbach, G. How sequential changes in reward expectation  
750 modulate cognitive control: Pupillometry as a tool to monitor dynamic changes in reward  
751 expectation. *Int. J. Psychophysiol.* 148, 35–49 (2020).
- 752 36. Yang, H., Bari, B. A., Cohen, J. Y. & O'Connor, D. H. Locus coeruleus spiking differently  
753 correlates with S1 cortex activity and pupil diameter in a tactile detection task. *Elife* 10,  
754 (2021).
- 755 37. Vinck, M., Batista-Brito, R., Knoblich, U. & Cardin, J. A. Arousal and locomotion make  
756 distinct contributions to cortical activity patterns and visual encoding. *Neuron* 86, 740–754  
757 (2015).
- 758 38. Mineault, P. J., Tring, E., Trachtenberg, J. T. & Ringach, D. L. Enhanced spatial resolution  
759 during locomotion and heightened attention in mouse primary visual cortex. *J. Neurosci.* 36,  
760 6382–6392 (2016).
- 761 39. McGinley, M. J., David, S. V. & McCormick, D. A. Cortical membrane potential signature of  
762 optimal states for sensory signal detection. *Neuron* 87, 179–192 (2015).
- 763 40. Stringer, C. et al. Spontaneous behaviors drive multidimensional, brainwide activity.  
764 *Science* 364, 255 (2019).
- 765 41. Lipton, D. M., Gonzales, B. J. & Citri, A. Dorsal striatal circuits for habits, compulsions and  
766 addictions. *Front. Syst. Neurosci.* 13, 28 (2019).
- 767 42. Mendelsohn, A. I. Creatures of habit: the neuroscience of habit and purposeful behavior.  
768 *Biol. Psychiatry* 85, e49–e51 (2019).
- 769 43. Amaya, K. A. & Smith, K. S. Neurobiology of habit formation. *Curr. Opin. Behav. Sci.* 20,  
770 145–152 (2018).
- 771 44. de Wit, S. & Dickinson, A. Associative theories of goal-directed behaviour: a case for  
772 animal-human translational models. *Psychol. Res.* 73, 463–476 (2009).
- 773 45. Smith, K. S. & Graybiel, A. M. A dual operator view of habitual behavior reflecting cortical  
774 and striatal dynamics. *Neuron* 79, 361–374 (2013).
- 775 46. Adams, C. D. Variations in the sensitivity of instrumental responding to reinforcer  
776 devaluation. *The Quarterly Journal of Experimental Psychology Section B* 34, 77–98  
777 (1982).
- 778 47. Coutureau, E. & Killcross, S. Inactivation of the infralimbic prefrontal cortex reinstates goal-  
779 directed responding in overtrained rats. *Behav. Brain Res.* 146, 167–174 (2003).
- 780 48. O'Doherty, J. P., Cockburn, J. & Pauli, W. M. Learning, reward, and decision making. *Annu.*  
781 *Rev. Psychol.* 68, 73–100 (2017).
- 782 49. Miller, K. J., Shenhav, A. & Ludvig, E. A. Habits without values. *Psychol. Rev.* 126, 292–  
783 311 (2019).

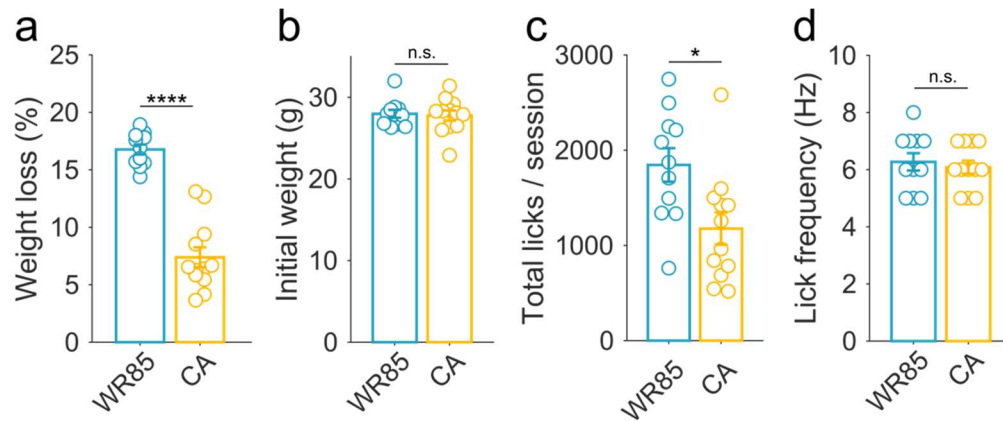
- 784 50. Watson, P. & de Wit, S. Current limits of experimental research into habits and future  
785 directions. *Curr. Opin. Behav. Sci.* 20, 33–39 (2018).
- 786 51. Garrett, N., Allan, S. & Daw, N. D. Model based control can give rise to devaluation  
787 insensitive choice. *BioRxiv* (2022). doi:10.1101/2022.08.21.504635
- 788 52. Vandaele, Y., Pribut, H. J. & Janak, P. H. Lever insertion as a salient stimulus promoting  
789 insensitivity to outcome devaluation. *Front Integr Neurosci* 11, 23 (2017).
- 790 53. Garr, E. & Delamater, A. R. Exploring the relationship between actions, habits, and  
791 automaticity in an action sequence task. *Learn. Mem.* 26, 128–132 (2019).
- 792 54. Thraillkill, E. A. & Bouton, M. E. Contextual control of instrumental actions and habits. *J.*  
793 *Exp. Psychol. Anim. Learn. Cogn.* 41, 69–80 (2015).
- 794 55. Dickinson, A., Nicholas, D. J. & Adams, C. D. The effect of the instrumental training  
795 contingency on susceptibility to reinforcer devaluation. *The Quarterly Journal of*  
796 *Experimental Psychology Section B* 35, 35–51 (1983).
- 797 56. Shillinglaw, J. E., Everitt, I. K. & Robinson, D. L. Assessing behavioral control across  
798 reinforcer solutions on a fixed-ratio schedule of reinforcement in rats. *Alcohol* 48, 337–344  
799 (2014).
- 800 57. Vandaele, Y. et al. Distinct recruitment of dorsomedial and dorsolateral striatum erodes with  
801 extended training. *Elife* 8, (2019).
- 802 58. Smith, K. S. & Graybiel, A. M. Using optogenetics to study habits. *Brain Res.* 1511, 102–  
803 114 (2013).
- 804 59. Smith, K. S. & Graybiel, A. M. Habit formation coincides with shifts in reinforcement  
805 representations in the sensorimotor striatum. *J. Neurophysiol.* 115, 1487–1498 (2016).
- 806 60. Wu, J. et al. Transcardiac perfusion of the mouse for brain tissue dissection and fixation.  
807 *Bio Protoc* 11, e3988 (2021).
- 808 61. Juczewski, K., Koussa, J. A., Kesner, A. J., Lee, J. O. & Lovinger, D. M. Stress and  
809 behavioral correlates in the head-fixed method: stress measurements, habituation  
810 dynamics, locomotion, and motor-skill learning in mice. *Sci. Rep.* 10, 12245 (2020).
- 811 62. Kuchibhotla, K. V. et al. Dissociating task acquisition from expression during learning  
812 reveals latent knowledge. *Nat. Commun.* 10, 2151 (2019).
- 813 63. Moore, S. & Kuchibhotla, K. V. Slow or sudden: Re-interpreting the learning curve for  
814 modern systems neuroscience. *IBRO Neuroscience Reports* 13, 9–14 (2022).
- 815 64. Nath, T. et al. Using DeepLabCut for 3D markerless pose estimation across species and  
816 behaviors. *Nat. Protoc.* 14, 2152–2176 (2019).
- 817 65. Mathis, A. et al. DeepLabCut: markerless pose estimation of user-defined body parts with  
818 deep learning. *Nat. Neurosci.* 21, 1281–1289 (2018).

819

820

821 **Extended Materials**

822 **Extended Figures**

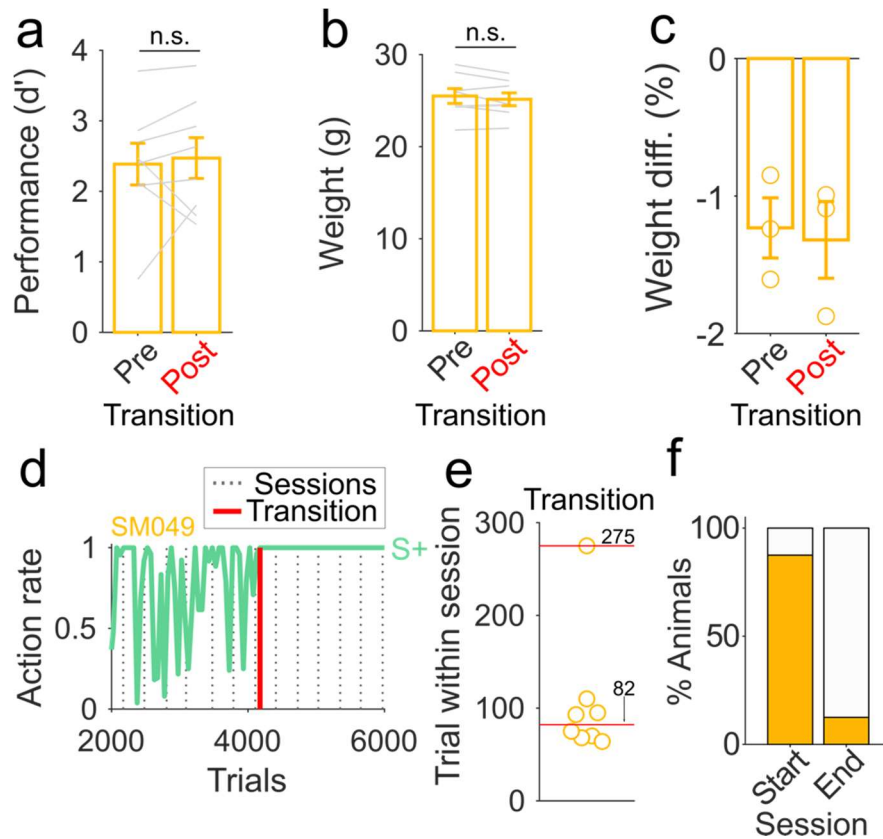


823

824 **Extended Data Figure 1. Palatability-based motivation is stable throughout training and reduces an**  
825 **animals' motivation to obtain water rewards.**

826 **a**, Weight loss during go/no-go training was significantly reduced in CA mice (Wilcoxon ranksum test,  
827  $p=0.000055$ ). **b**, Animals' original weights were not different between groups (Wilcoxon ranksum test,  
828  $p=0.97$ ). **c**, CA mice also do significantly lower number of total licks in a lick-training session (Wilcoxon  
829 rank-sum test,  $p=0.021$ ). **d**, Lick vigor (frequency) was not different between groups ( $p=0.67$ ). **a-d**,  $n=11$   
830 WR85 mice and  $n=12$  CA mice.

831

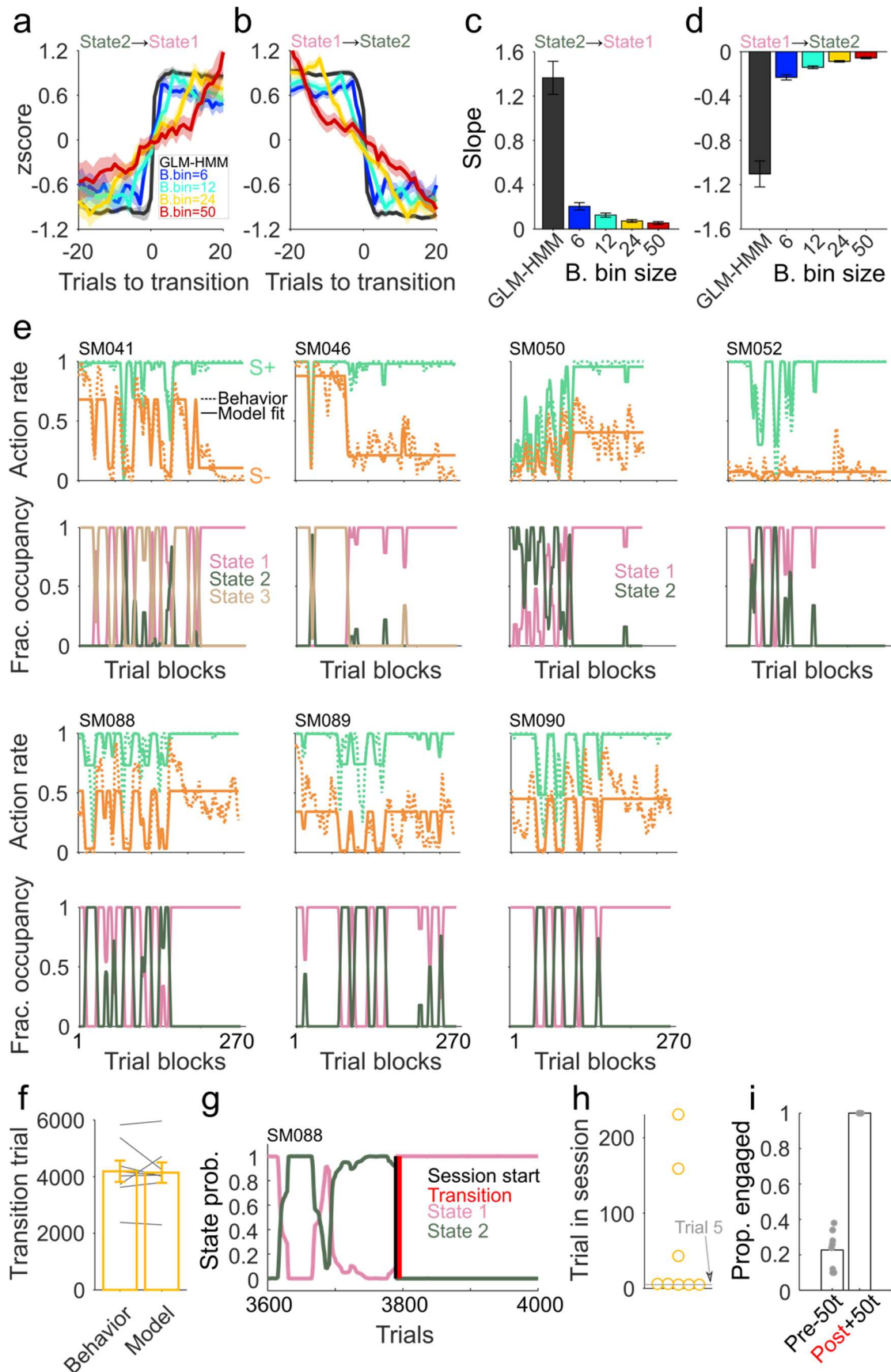


832  
 833  
 834  
 835  
 836  
 837  
 838  
 839  
 840  
 841  
 842  
 843  
 844  
 845  
 846  
 847

**Extended Data Figure 2. Spontaneous transitions from goal-directed to habitual behavior are independent of performance or metabolic state and occur at the beginning of a new session.**

**a**, No significant differences are observed between pre and post-transition performance (Wilcoxon signed rank test,  $p=0.54$ ,  $n=8$  mice). **b**, Weights of CA mice with transitions are stable around the transition session (average of 3 sessions pre and post transition). No significant differences are observed between pre-transition weight and post transition weight (Wilcoxon signed rank test,  $p=0.29$ ,  $n=8$  mice). **c**, Weight differences between the end of a session and before the start of the next day session confirm that animals maintain similar consumption rates in their home-cage comparing pre and post-transition (Wilcoxon signed rank test,  $p=0.82$ ,  $n=3$  mice). **d**, Exemplar animal of hit rates (green) with the overlying sessions (vertical dotted gray lines) showing that the transition to habit happened close to the start of a new session. **e**, Of the 8 CA mice with transitions, 7 transitioned early in the session. These data use a 50-trial binning procedure that limits the temporal resolution, such that animals transitioning at Trial ~80, are likely transitioning much closer to the beginning of the session. See Extended Data Figure 3g-h for a model-based estimation using trial-level data. **f**, The vast proportion of CA mice transition at the start of a new session (trials 1-150), compared to in the end of a session (trial 151-300) ( $n=8$  mice).

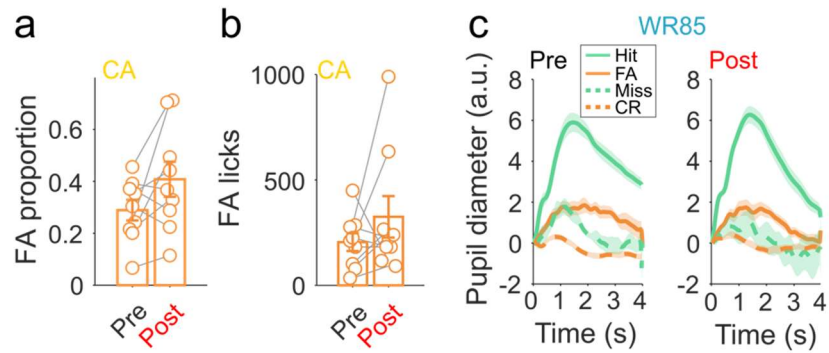






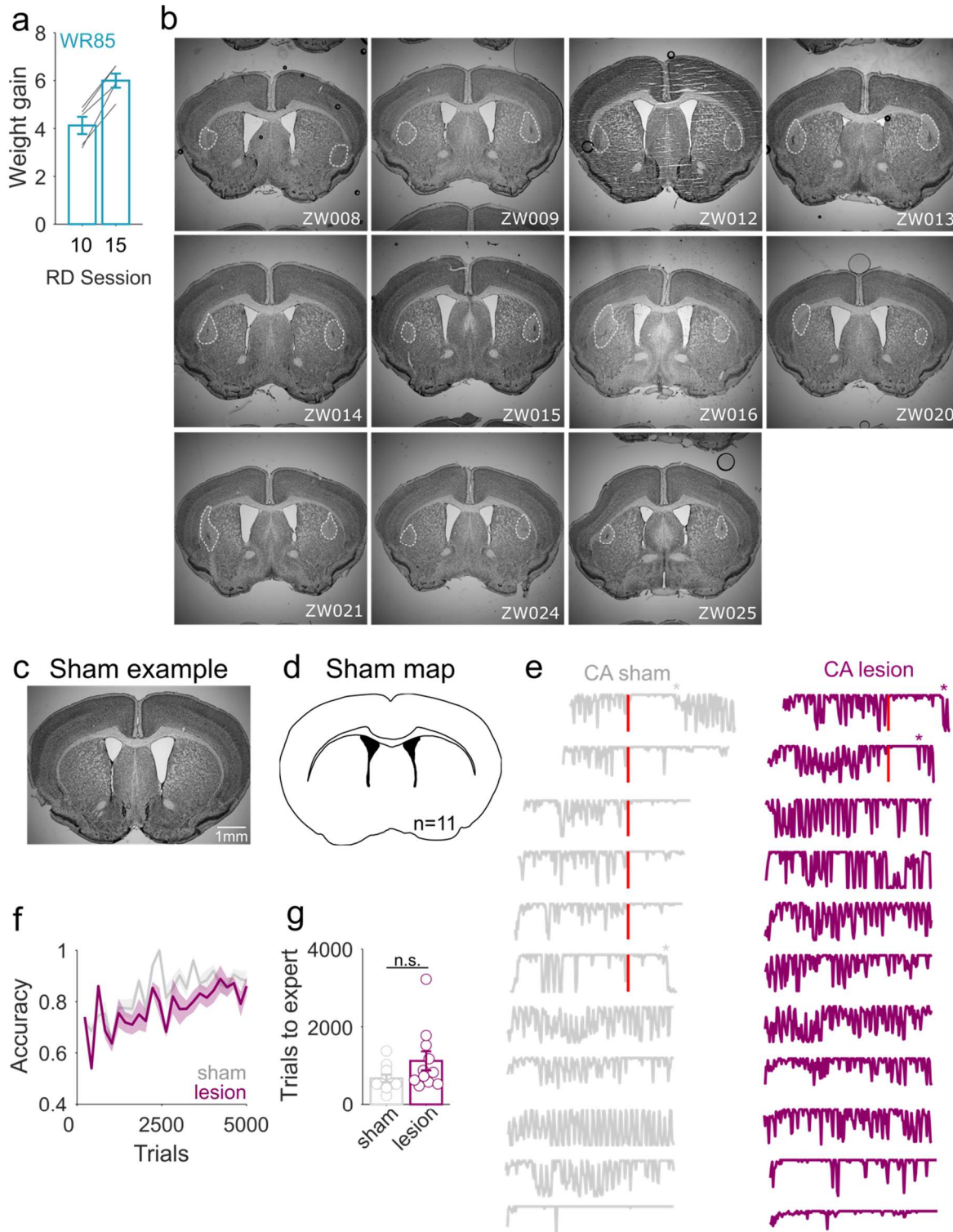
849 **Extended Data Figure 3. State-like transitions are observed in individual animals using an HMM-**  
850 **GLM model.**

851 Z-scored model state probability (black) or behavioral hit rate using bin size from 6-50 trials (blue – red) for  
852 **a**, disengaged-to-engaged and **b**, engaged-to-disengaged fluctuations. Quantification of abruptness using  
853 slope of the trajectories near the transition point for **c**, disengaged-to-engaged and **d**, engaged-to-  
854 disengaged fluctuations. **e**, Model fit plots and fraction occupancy plots for all additional CA individual  
855 animals that transitioned (n=7). **f**, The trial-by-trial nature of the model allowed us to find exact transition  
856 points, which are similar to the behaviorally identified ones using only hit rates (Wilcoxon signed rank test  
857  $p=0.95$ ). **g**, State probability for State 1 (pink), State 2 (green) for an example animal right around the  
858 transition session (black vertical line), depicting that transitions (red vertical line) typically happen right at  
859 the start of a new session. **h**, Most animals (5/8) transition within 5 trials, but others transition later in a  
860 session. **i**, 50 trials before the transition point (Pre -50t) there is very low probability of task engagement,  
861 while 50 trials after the transition (Post +50t), engagement is at its maximum.  
862  
863



864

865 **Extended Data Figure 4. Changes in pupillary error signal between pre and post transition are only**  
866 **evident in CA mice and independent of movement-evoked pupillary changes.** **a**, The proportion of FA  
867 in CA mice is not different pre and post-transition (Wilcoxon signed rank test,  $p=0.098$ ) ( $n=3$  days pre and  
868 3 days post from 4 mice). **b**, No changes in the number of licks to FA were seen in CA mice pre and post-  
869 transition (Wilcoxon signed rank test,  $p=0.25$ ) ( $n=3$  days pre and 3 days post from 3 mice). **c**, Tone-evoked  
870 pupil dilation during hits or FA is not changed between pre and post-transition in WR85 mice ( $n=3$  days pre  
871 and 3 days post from 4 mice).



873 **Extended Data Figure 5. DLS lesioned animals are less likely to transition to habitual behavior.**  
874 **a**, No decrease in weight between reward devaluation (RD) session 10 and 15, suggesting that differences  
875 in the increased hit rates seen in session 15 are not due to reduced water consumption during the satiety  
876 test (Wilcoxon signed rank test,  $p=0.0039$ ) ( $n=5$  WR85 mice). **b**, CA-sham exemplar showing no DLS  
877 lesions when injected bilaterally with vehicle. **c**, Map of  $n=11$  CA-sham mice showing no lesion-like areas  
878 for any individual. **d**, All individual lesioned mice with outlined lesion area (dotted white line). **e**, Sham (gray)  
879 and lesioned (purple) individual mouse hit rates, depicting a transition to habitual behavior (red). Animals  
880 that transition back to goal-directed control are shown with an asterisk. **f**, No differences in task accuracy  
881 between sham (gray) and lesioned (purple) animals (2-way ANOVA,  $F(1,25)=0.43$ ,  $p=0.51$ , interaction  
882 group  $\times$  trials  $F(1,25)=0.5$ ,  $p=0.98$ ) ( $n=11$  sham and  $n=11$  lesioned mice). **g**, Similar number of trials to  
883 expert performance between sham (gray) and lesioned (purple) mice (Wilcoxon ranksum test  $p=0.098$ )  
884 ( $n=11$  sham and  $n=11$  lesioned mice).