

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

Prelimbic cortex maintains attention to category-relevant information and flexibly updates category representations

Matthew B. Broschard¹, Jangjin Kim¹, Bradley C. Love², Edward A. Wasserman¹, John H. Freeman¹

¹Department of Psychological and Brain Sciences, University of Iowa, Iowa City, IA, USA, 52242

²Department of Experimental Psychology and The Alan Turing Institute, University College London, London, UK,

Correspondence

Matthew B. Broschard, Department of Psychological and Brain Sciences, University of Iowa, Iowa City, IA 52242. Email: matthew-broschard@uiowa.edu

Funding information

National Institutes of Health Grant P01-HD080679 to J.H.F. and B.C.L. and Wellcome Trust Investigator Award WT106931MA to B.C.L.

Keywords: prelimbic prefrontal cortex, rat, category learning, executive functions, SUSTAIN, touchscreen

42 **Highlights**

- 43 • Rats categorize distribution of stimuli containing two continuous dimensions
- 44 • Prefrontal lesions impair category tasks containing irrelevant stimulus information
- 45 • Prefrontal lesions do not affect tasks containing only relevant stimulus information
- 46 • Prefrontal lesions impair trial-by-trial updating of category representations

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65 **Abstract**

66 Category learning groups stimuli according to similarity or function. This involves finding and
67 attending to stimulus features that reliably inform category membership. Although many of the
68 neural mechanisms underlying categorization remain elusive, models of human category learning
69 posit that prefrontal cortex plays a substantial role. Here, we investigated the role of the
70 prelimbic cortex (PL) in rat visual category learning by administering excitotoxic lesions before
71 category training and then evaluating the effects of the lesions with computational modeling.
72 Using a touchscreen apparatus, rats (female and male) learned to categorize distributions of
73 category stimuli that varied along two continuous dimensions. For some rats, categorizing the
74 stimuli encouraged selective attention towards a single stimulus dimension (i.e., 1D tasks). For
75 other rats, categorizing the stimuli required divided attention towards both stimulus dimensions
76 (i.e., 2D tasks). Testing sessions then examined generalization to novel exemplars. PL lesions
77 impaired learning and generalization for the 1D tasks, but not the 2D tasks. Then, a neural
78 network was fit to the behavioral data to examine how the lesions affected categorization. The
79 results suggest that the PL facilitates category learning by maintaining attention to category-
80 relevant information and updating category representations.

81

82

83 Categorization is the process of grouping perceptually or functionally related objects and events.
84 Abundant evidence from neuroimaging (Kumaran, Summerfield, Hassabis, & Maguire, 2009;
85 Bowman & Zeithamova, 2018) and physiology (Freedman et al., 2001) experiments supports
86 the recruitment of prefrontal cortex (PFC) in categorization tasks. The PFC is also important for
87 transitive inference, a mechanism that infers new information and promotes generalization by

88 extrapolating overlapping information across multiple episodes (Koscik & Tranel, 2012;
89 Zeithamova, Dominick, & Preston, 2012).

90 Accordingly, theories of categorization predict that the PFC plays a substantial role in
91 learning new categories. COVIS (COmpetition between Verbal and Implicit Systems) posits that
92 the PFC governs a declarative system that learns new categories by testing explicit category rules
93 (Ashby et al., 1998). The COVIS framework has been tested empirically by training participants
94 to categorize distributions of visual stimuli that vary along two continuous dimensions (Maddox,
95 Ashby, & Bohil, 2003; Smith et al., 2012). In one condition, only one stimulus dimension is
96 category-relevant, and learning involves *selective attention* to that dimension (1D tasks; Fig. 1B).
97 In a second condition, both stimulus dimensions are relevant, and learning requires *divided*
98 *attention* to both dimensions (2D tasks; Fig. 1C). COVIS predicts that the declarative system
99 (and the PFC) is important for learning 1D tasks, as they can be solved by a unidimensional
100 category rule (Ashby & Maddox, 2011). This prediction is supported by neuroimaging
101 experiments (Nomura et al., 2006).

102 Rodents have become great models to examine mechanisms underlying complex
103 behavior (Zoccolan, Oertelt, DiCarlo & Cox, 2009; Vinken, Vermaercke & Op de Beeck, 2014).
104 We recently developed rodent versions of the 1D and 2D tasks using a touchscreen apparatus to
105 investigate rat category learning (Broschard, Kim, Love, Wasserman, & Freeman, 2019). The
106 current experiment extends this work by examining the contributions of the prelimbic (PL) area
107 of the rat PFC. Broschard et al., 2019 concluded that rats use selective attention to learn the 1D
108 tasks and bias attention towards the category-relevant dimension. We predict that this is
109 mediated by the PL; therefore, inactivating the PL will impair learning for the 1D tasks. This
110 prediction is supported by calcium imaging in the mouse medial frontal cortex during a go/no-go

111 version of the 1D task (Reinert et al., 2021). This prediction also aligns with Love & Gureckis
112 (2007), who proposed that the PFC is synonymous to the selective attention mechanism of the
113 neural network model SUSTAIN (Supervised and Unsupervised STRatified Adaptive Incremental
114 Network; Love, Medin, & Gureckis, 2004). The current experiment tested this prediction
115 directly.

116 There is contention regarding whether rodent PL is comparable to the primate PFC
117 (Laubach, Amarante, Swanson, & White, 2018). PL satisfies early definitions of PFC by
118 exhibiting bidirectional communication with the medial dorsal thalamus (Rose & Woolsey,
119 1948). Additionally, some functions of PL are analogous to primate PFC, including working
120 memory (Horst & Laubach, 2009), goal directed behavior (Ostlund 2005), response conflict
121 (Wit, Kosaki, Balleine, & Dickinson, 2006), behavioral flexibility (Ragozzino 2007), and
122 attention (Tait, Bowman, Neuwirth & Brown, 2018). However, anatomical investigations
123 conclude that PL may be homologous to cingulate cortex in primates (Heilbronner et al., 2016).
124 Furthermore, all of rodent frontal cortex is agranular, highlighting large differences in the
125 cellular makeup between rodents and primates (Uylings & Eden, 1991; Seamans, Lapish &
126 Durstewitz, 2008). Therefore, generalizing the results of the current experiment to primate PFC
127 requires careful consideration of anatomical and functional comparisons.

128 Here, we investigated the role of the PL in visual category learning in rats. Rats
129 underwent stereotaxic surgery to lesion the PL with NMDA. After recovery, the rats were trained
130 to learn the 1D or 2D categorization tasks. Then, we fit the neural network SUSTAIN to the
131 behavioral data to further examine the role of the PL, specifically as it pertains to selective
132 attention. Together, the results suggest that the PL maintains attention to category-relevant
133 information and updates category representations according to recent exemplars.

134

135 **Materials and Methods**

136 *Subjects*

137 Male (n = 16, mean weight: ~350 grams) and female (n = 16, mean weight: ~250 grams) Long-
138 Evans rats were studied. Upon arriving in the animal colony, rats were put on a 12-hour
139 light/dark cycle and given *ad libitum* access to food and water. After acclimating to the new
140 environment for a week, food was restricted. Weights were recorded daily to ensure weights did
141 not go below 85% of the rats' free feeding weight. All procedures were approved by the
142 Institutional Animal Care and Use Committee at the University of Iowa.

143

144 *Touchscreen Apparatus*

145 For all experimental sessions, rats were placed within custom-built touchscreen chambers
146 (Figure 1A; 36 × 41 × 36 cm). The chambers contained a computer monitor (Model 1550V,
147 NEC, Melville, NY) mounted on one wall to present visual stimuli to the rats. A touchscreen
148 (15-in, Elo Touch Systems, Fremont, CA) was placed in front of the computer monitor so that
149 the rats could interact with the screen. On the wall opposite from the monitor, a food tray (6.5 ×
150 13 × 4.5 cm) delivered food pellets to the rat via a rotary pellet dispenser (Med Associates Inc.,
151 Georgia, VT, model ENV-203IR) that was controlled by an electrical board (Model RS-232,
152 National Control Devices, Osceola, MO). A house light above the food tray was always on
153 during experimental sessions. White noise within the room was also always on to minimize
154 distractions. Custom MATLAB scripts controlled all experimental sessions and procedures
155 (MathWorks, Natick, MA). Finally, a camera (model ELP-USB100W05MT-RL36) was mounted

156 to the ceiling of the chamber and faced the computer screen so that the rats' behavior could be
157 observed and recorded.

158

159 *Pre-Training Procedures*

160 Once food restriction began, each rat was handled daily for 1 week. This reduced the stress of
161 interacting with experimenters. Then, each rat underwent cart training, which encouraged the
162 foraging of food pellets in an open field. Each rat was placed on the surface of a laboratory cart,
163 and twenty 45-mg pellets were scattered on the cart's surface. This procedure was repeated daily
164 until the rat consumed all pellets within 15 minutes, which usually took about 7 days. After cart
165 training, rats underwent a daily shaping procedure to learn to interact with the touchscreen
166 (Broschard, Kim, Love, & Freeman, 2020). This procedure included three separate phases; each
167 phase was incrementally similar to the trial sequence used during training and testing sessions.
168 Phase I required a minimal touch requirement and was used to orient the rats to the screen. Each
169 trial began with the presentation of a star at the center of the screen. After 15 seconds (or one
170 touch of the screen), the star was replaced by a white box appearing on the left or right side of
171 the screen. A food pellet was delivered if the rat touched anywhere on the screen while the white
172 box was presented. Otherwise, the trial aborted after 45 seconds, and the trial was considered a
173 miss. This procedure was repeated until the rat completed at least 55/60 trials within 25 minutes.
174 In Phase II, the touch requirement was increased. Specifically, the rats were required to touch
175 both the star stimulus and the white box to receive a food reward. Similar to Phase I, the trial
176 phases timed out (i.e., 15 seconds for the star stimulus and 45 seconds for the white box) in the
177 absence of a response. Sessions continued until the rat completed at least 55 trials within 30
178 minutes. Phase III was identical to Phase II except that the trials did not time out. Sessions

179 continued until the rat completed all 60 trials within 25 minutes. All shaping procedures required
180 about 14 days.

181

182 *Surgery*

183 After shaping was complete, rats underwent stereotaxic surgery. Under isoflourane (1% - 4%)
184 anesthesia), a Hamilton syringe (1 μ L; 26 gauge) was lowered into the PL bilaterally (AP: +3.0;
185 ML: \pm 0.7; DV; -3.5). Upon reaching the target site, 0.4 μ L of either NMDA (20 mg/ml; 10 μ L/h;
186 Sigma-Aldrich, St. Louis, MO) or PBS was infused. After surgery, rats were placed on a heating
187 pad until awake and mobile to prevent hypothermia. Meloxicam (1 mg/ml) was administered as
188 analgesic both during surgery and 24 hours after surgery. Rats were allowed at least one week to
189 recover.

190

191 *Behavioral Testing: An Overview*

192 After a week of recovery, rats were given multiple training and testing sessions to learn to
193 categorize visual stimuli. Briefly, on each trial, a single stimulus appeared on the screen, and the
194 rat decided its category membership (i.e., category 'A' or category 'B') by pressing one of two
195 report keys (Fig. 1D). Food reinforcement was delivered after correct responses to guide
196 learning.

197

198 *Category Stimuli*

199 The category stimuli (239 x 239 pixels) presented to the rats contained black and white gratings
200 (Figs. 1B-D). Across stimuli, these gratings varied along two continuous dimensions: spatial
201 frequency and orientation. The spatial frequency of the gratings ranged from 0.2532 cycles per

202 visual degree (cpd) to 1.2232 cpd, and the orientation of the gratings ranged from 0 radians to
203 1.75 radians. These values were obtained from pilot experiments and are within the perceptual
204 limits of rats (Crijns & Op de Beeck, 2019). Linear transformations of these dimensions were
205 made so that both dimensions had a common range (i.e., 0 to 100). Specifically,

206
$$\text{Normalized frequency} = \frac{\text{cpd}}{0.0097} - 26.10,$$

207
$$\text{Normalized orientation} = \text{radians} * \frac{180}{\pi}.$$

208 A two-dimensional stimulus space was created using these transformed stimulus dimensions
209 (Figs. 1B-C).

210

211 *Category Tasks*

212 Category tasks were created by placing bivariate normal distributions on this transformed
213 stimulus space (Fig. 1B; Category A: $\mu_X = 30$, $\sigma_X = 2.5$, $\mu_Y = 50$, $\sigma_Y = 20$; Category B: $\mu_X = 70$, σ_X
214 $= 2.5$, $\mu_Y = 50$, $\sigma_Y = 20$; Broschard et al., 2019; Broschard et al., 2020; O'Donoghue, Broschard,
215 & Wasserman, 2020). Each distribution constituted a category, and each point within a
216 distribution represented a category stimulus. Three additional category tasks were created by
217 rotating these distributions in 45-degree increments (Figs. 1B-C). Importantly, rotating the
218 distributions did not affect any physical properties of the distributions (Ashby, Smith, &
219 Rosedahl, 2019; e.g., standard deviation, mean between-category distance, etc.). However, these
220 rotations changed how the distributions were oriented in relation to the axes of the stimulus
221 space. 1D tasks had distributions that were perpendicular to one of the stimulus dimensions (Fig.
222 1B). Because of this orientation, only one dimension (i.e., the perpendicular dimension) was
223 category-relevant and had to be considered when deciding category membership. The dimension
224 parallel to the distributions was category-irrelevant and could be ignored. Conversely, 2D tasks

225 had distributions that were not aligned with either stimulus axis (Fig. 1C). For these tasks, both
226 dimensions were category-relevant, and deciding category membership involved combining
227 information from both dimensions.

228

229 *Category Training*

230 Rats were randomly assigned to learn one of the four category tasks (Broschard et al., 2019,
231 2020). Rats were given 15 training sessions; each session contained 80 training trials. On each
232 trial, a star stimulus was presented at the center of the screen (Fig. 1D; Star Phase). After one
233 touch of the star, a category exemplar was randomly selected from the training distributions
234 (Figs. 1B-C) and replaced the star stimulus (Cue Phase). After three touches of this exemplar,
235 copies of the exemplar were presented on the left and right sides of the screen, acting as report
236 keys (Choice Phase). Rats touched either report key depending on the category membership of
237 the exemplar. The categories were mapped spatially, such that the left report key was chosen for
238 members of category A, and the right report key was chosen for members of category B. If the
239 correct side was chosen, a white box replaced the report key (Reward Phase). One touch of the
240 white box delivered a food reward. If instead the incorrect side was chosen, then a correction
241 trial was initiated. Here, the trial repeated from the Cue Phase after a 5 to 10 second time-out.
242 Correction trials were repeated without reinforcement until the correct side was chosen. Inter-
243 trial intervals ranged from 5 to 10 seconds.

244

245 *Category Generalization*

246 After category training, rats were presented with five testing sessions to examine category
247 generalization (Broschard et al., 2019, 2020). Each session contained 80 trials. The trial sequence

248 was identical to training sessions except that correction trials were not administered after
249 incorrect responses (and therefore all choices were reinforced). Exemplars were randomly
250 sampled from testing distributions (Fig. 5A). Testing distributions were identical to the training
251 distributions, except that the standard deviation along the relevant dimension (or axis for the 2D
252 tasks) was increased ($\sigma_X = 10$; Broschard et al., 2019; O'Donoghue et al., 2020). With this
253 manipulation, some exemplars overlapped with the training distributions (i.e., Trained; within
254 two standard deviations), but some exemplars sampled from novel portions of the stimulus
255 space. Among the novel exemplars, about half were closer to the category boundary than the
256 training distributions (Proximal), and half were farther from the category boundary (Distal).
257 Generalization to the novel stimuli ensures that the rats did not simply memorize single
258 exemplars during training.

259

260 *Simple Discrimination*

261 After category testing, rats were trained to learn a simple discrimination task. This acted as a
262 control task to ensure that any differences across groups were not caused by deficits in
263 movement, motivation, perception, etc. Instead of categories of stimuli, only two images were
264 presented during training sessions (i.e., a light box and a dark box; Fig. 6A; Kim, Castro,
265 Wasserman, & Freeman, 2018). Both images contained a common pattern of dots to add
266 perceptual complexity. The trial sequence was identical to categorization sessions. The white
267 stimulus was mapped to the left report key, and the black stimulus was mapped to the right report
268 key. Each session contained 72 training trials. Sessions continued until the rat reached a learning
269 criterion (i.e., at least 75% accuracy for both images on two consecutive sessions).

270

271 *Histology*

272 After all behavioral testing, rats were perfused to verify lesion placements. Rats were given a
273 lethal dose of euthanasia solution (sodium pentobarbital) and then perfused with ~400 mL PBS
274 and ~400 mL of 10% formalin. Brains were stored at 4° C in a solution containing 10% formalin
275 and 30% sucrose. A sliding microtome collected 50 µm coronal sections of the target area. Brain
276 sections were then stained with thionin (Sigma-Aldrich, St. Louis, MO). A close investigation of
277 the tissue was conducted under a light microscope to characterize the size of each lesion within
278 the PL and how much it extended dorsally and ventrally. The boundary of the PL was defined
279 according to Paxinos & Watson, 1998.

280

281 *Statistical Analysis*

282 Multiple dependent measures quantified performance for training and testing sessions. First,
283 session accuracy was defined as the proportion of correct responses during the Choice phase.
284 Second, perseverative errors were calculated and were defined as a repeated incorrect response
285 after receiving negative feedback. Third, reaction time was calculated during the Cue phase and
286 Choice phase to quantify the amount of time to 1) observe the stimulus and 2) make a category
287 decision. Reaction times from incorrect trials were excluded from all analyses. Additionally,
288 reaction times that exceeded two standard deviations of the mean were excluded from all
289 analyses, a criterion that is commonly used to eliminate outliers (O'Donoghue et al., 2020).
290 These outliers rarely occurred. Fourth, touch separation used the pixel location of touches during
291 the Cue phase of correct trials to quantify choice confidence. Prior experiments demonstrated
292 that as accuracy improves, the x-coordinate of touches during the Cue phase deviate towards the
293 correct side in anticipation of the rats' choice (Kim, Castro, Wasserman, & Freeman, 2018).

294 Touch separation is calculated by comparing the x-coordinate of a touch to the average x-
295 coordinate of all three touches from that trial. Positive touch separation indicates deviation
296 towards the correct side, and negative touch separation indicates deviation towards the incorrect
297 side.

298 These dependent measures were analyzed using linear mixed effects modeling (R,
299 version 3.4.2). Models used for training sessions included fixed effects for experimental group,
300 training session, and a quadratic function across training sessions, as well as random effects for
301 slope, intercept, and the quadratic function. Models for testing sessions included fixed effects for
302 experimental group, trial type (Distal, Trained, and Proximal), and a quadratic function across
303 trial types, as well as random effects for slope, intercept, and the quadratic function. Quadratic
304 functions were used because they best fit the data, and higher order terms did not significantly
305 improve these fits. Sex was added as a covariate for all models to check whether there were any
306 significant differences between male and female rats. To find the simplest model that fit the data,
307 we used a model simplification strategy (Crawley 2007). We started with the full model and then
308 systematically removed random effects one at a time. This continued until the estimates were
309 significantly different from the larger model before it.

310

311 *SUSTAIN Model Fitting*

312 SUSTAIN is a neural network model of human category learning and has been used in multiple
313 contexts to map neural activity to specific cognitive processes (e.g., Love & Gureckis, 2007;
314 Mack, Love & Preston, 2016). Here, we used SUSTAIN to further examine the role of the PL by
315 simulating the effects of the PL lesions on category learning. We were particularly interested

316 whether the PL serves a function similar to SUSTAIN's attention mechanism (Love & Gureckis,
317 2007).

318 SUSTAIN assumes that similar training experiences cluster together in memory (Love et
319 al., 2004). Categories are represented by one or multiple clusters; each cluster reflects a learned
320 group of similar training experiences and is stored in a hidden layer (Fig. 7A; the cluster layer).
321 On each learning trial, the current stimulus is compared to existing clusters, and each cluster is
322 activated according to its similarity to the stimulus. SUSTAIN's attention mechanism modulates
323 the stimulus before entering the cluster layer (Fig. 7A; the feature tuning mechanism). Each
324 stimulus dimension is multiplied by an attention weight. These weights bias the perception of the
325 stimulus according to category-relevant information and affect how clusters are activated.
326 Cluster activations then project to a decision layer, which makes a probabilistic decision
327 regarding the category membership of the stimulus (Fig. 7A; decision layer).

328 At the beginning of training, the model contains one cluster centered on the first training
329 stimulus, and attention weights are equivalent across all stimulus dimensions. Then, feedback is
330 provided after each trial, and SUSTAIN updates accordingly. First, category representations
331 within the cluster layer update, such that the current trial stimulus is either integrated into an
332 existing cluster or becomes the center of a newly recruited cluster. New clusters are created in
333 response to stimuli that are 'surprising.' The decision to recruit a new cluster is initiated if the
334 model incorrectly classifies a stimulus and the cluster activations exceed the value of a threshold
335 parameter, indicating that the model is relatively confident in its choice. The feature tuning
336 mechanism is also updated so that attention is shifted towards category-relevant dimensions.
337 This is controlled by two parameters. First, a selective attention parameter determines the
338 amount of attentional focus that can be applied in the category task. Second, an attention learning

339 rate parameter determines how quickly this attention resource can be shifted towards relevant
340 dimensions.

341 Love & Gureckis (2007) proposed a framework by which the functions of the PFC map
342 onto elements of the SUSTAIN model. Specifically, they posit that the PFC functions as the
343 feature tuning mechanism and shifts attention towards category-relevant information. Second,
344 the PFC updates category representations by initiating the decision to recruit a new cluster. To
345 test these predictions, we created three experimental manipulations that simulate the effects of
346 the PL lesions. The first two manipulations disrupted the feature tuning mechanism to test
347 whether the PL is critical for shifting attention to relevant dimensions. First, we lesioned the
348 feature tuning mechanism by setting the two parameters that control the feature tuning
349 mechanism (i.e., the selective attention parameter and the attention learning rate parameter) to 0.
350 As a result, the model could not update its attention weights, rendering the model unable to shift
351 attention to category-relevant dimensions. Second, we permuted the attention weights before
352 each trial. With this manipulation, the model could update its attention weights normally;
353 however, on any given trial, attention may be directed towards category-irrelevant information.
354 Therefore, the model could learn to identify relevant information, but its ability to maintain
355 selective attention to that information across trials was impaired. The third manipulation tested
356 the prediction that the PL initiates the decision to recruit a new cluster in response to ‘surprising’

357 stimuli. This was accomplished by increasing the cluster threshold parameter that determines
358 when a new cluster is recruited.

359 Using combinations of these manipulations, we generated five versions of SUSTAIN that
360 each simulated how the PL lesions affected category learning (Fig. 7B). We also added a control
361 model that assumed the lesions had no effect on learning. Each model was optimized to the rats'
362 averaged learning curves using the MATLAB function *fmincon*. Then, Akaike's Information
363 Criterion (AIC) was calculated for each optimized model to quantify its goodness-of-fit (Akaike,
364 1974). The model with the smallest AIC value was determined as the model that best fit the
365 behavior. The function(s) of PL can be inferred from these results.

366

367 *Perceptual Recency Effect*

368 With the current design, each rat completed a large number of training trials. This afforded us the
369 ability to examine category learning on a trial-by-trial basis. Importantly, this sensitivity was
370 leveraged to further test the prediction that the PFC updates category representations (Love &
371 Gureckis, 2007). We examined the effect of the PL lesions on perceptual recency effects, which
372 characterize how category performance is influenced by the identity of the most recent training
373 exemplar (Jones, Love, & Maddox, 2006). Recency effects suggest that category decisions are
374 biased towards recent exemplars, which would imply that the learner regularly updates category
375 representations. Assuming representational updating is a function of the PFC, we predicted that
376 recency effects are mediated by the PL.

377 Recency effects often interact with the perceptual similarity between exemplars. For
378 example, performance is facilitated if the exemplar is perceptually similar to the most recent
379 exemplar (Jones et al., 2006). Therefore, we binned the accuracy¹ of training trials according to

¹ Perceptual recency effects can also be calculated by examining repeated responses rather than trial accuracy. We choose trial accuracy to simplify the measure. Additionally, because there are only two categories in the current design, results look similar using either method.

380 the perceived similarity between the current exemplar (n) and the most recent exemplar ($n-1$;
381 Nosofsky, 1986). Perceptual similarity between exemplars i and j was calculated as:

$$382 \quad s_{ij} = e^{-d_{ij}},$$

383 where d is the psychological distance between exemplars i and j . Psychological distance was
384 defined as,

$$385 \quad d_{ij} = \sum_{m=1}^M w_m * |x_i - x_j|$$

386 where w_m was SUSTAIN's estimated attention weight for dimension m on trial n , and x was the
387 physical value of the exemplar along dimension m . Trial effects were isolated by subtracting the
388 binned accuracies by the average of 1,000 permutations where trial order was shuffled.

389 Therefore, positive recency scores indicate increased accuracy due to trial order, negative scores
390 indicate decreased accuracy due to trial order, and 0 indicates no effect of trial order.

391

392 **Results**

393 *Histological assessment of PL lesions*

394 Representative lesions are shown in Figure 2. Each lesion was examined under a light
395 microscope to ensure that it was contained within the PL. PL boundaries were determined
396 according to Paxinos & Watson (1998). All lesions were centered within the PL, and the data
397 from all rats were included in all analyses. Along the rostral/caudate axis, all lesions were
398 contained between bregma +4.3 and +2.4. There were no significant differences in lesion size
399 and location between the males and females. The lesions of three rats (one rat learning a 1D task
400 and two rats learning a 2D task) extended dorsally into the cingulate cortex and ventrally into the
401 infralimbic cortex. However, there were no differences in behavior between rats with these
402 lesions and rats with more selective lesions.

403

404 *PL lesions impair category learning for 1D tasks, but not 2D tasks*

405 All rats completed 15 training sessions to learn either a 1D task or a 2D task. We used linear
406 mixed effects models to examine accuracy, the number of correction trials, and the number of
407 perseverative errors across category training (see Materials & Methods). The full models
408 included fixed effects for group, training session, a quadratic function (across sessions), random
409 effects for the intercept, slope, and the quadratic function, and a covariate for sex. For all
410 measures, there was a significant main effect for training session (Fig. 3). Session accuracy
411 increased across training, and the number of correction trials and perseverative errors decreased
412 across training (Accuracy: $t(27.11) = 5.20, p < .001$; Correction trials: $t(27.04) = 5.81, p < .001$;
413 Perseverative errors: $t(27.27) = 5.12, p < .001$). There were no significant differences between
414 male and female rats (Accuracy: $t(22.01) = -1.64, p = .116$; Correction trials: $t(19.25) = 0.67, p =$
415 $.513$; Perseverative errors: $t(29.01) = 0.46, p = .649$), suggesting that sex did not affect category
416 learning. There were also no significant differences between controls learning the 1D tasks vs.
417 the 2D tasks (Accuracy: $t(26.55) = .05, p = .963$; Correction trials: $t(26.78) = 0.04, p = .971$;
418 Perseverative errors: $t(27.02) = 0.46, p = .647$). This replicates our previous work and suggests
419 that rats normally learn 1D tasks and 2D at the same rate (Broschard et al., 2019).

420 Compared to controls, rats with PL lesions were impaired in learning the 1D tasks.
421 Specifically, accuracy was impaired, and the number of correction trials and perseverative errors
422 were larger (Figs. 3A-F; Accuracy: $t(27.40) = 2.43, p = .022$; Correction trials: $t(27.33) = 2.31, p$
423 $= .028$; Perseverative errors: $t(27.54) = 2.56, p = .030$). Conversely, PL lesions did not affect
424 category learning for the 2D tasks (Accuracy: $t(27.01) = 0.62, p = .541$; Correction trials:
425 $t(26.94) = 0.21, p = .838$; Perseverative errors: $t(26.87) = 0.33, p = .742$). Together, these results

426 indicate that the PL lesions impaired category learning for the 1D tasks, but not the 2D tasks.
427 The 1D tasks, but not the 2D tasks, involve category-irrelevant information, and therefore
428 encourage a shift in attention to a single stimulus dimension. Therefore, our results suggest that
429 the PL is important for shifting attention towards category-relevant dimensions and away from
430 irrelevant dimensions (i.e., selective attention). Without the PL, attention may be divided
431 between the relevant and irrelevant dimensions. Under this interpretation, the PL lesions did not
432 affect learning the 2D tasks because, without the PL, rats were biased toward deploying the
433 optimal strategy (i.e., divided attention) as both dimensions were relevant.

434

435 *Rats with PL lesions learning 1D tasks require more time to categorize exemplars*

436 Next, we examined the amount of time to evaluate each stimulus (Cue RT) and to execute a
437 category decision (Choice RT) using linear mixed effects models (fixed effects: group, training
438 session, a quadratic function (across sessions); random effects: intercept, slope, and the quadratic
439 function; covariate: sex). There were significant main effects of training session for both Cue RT
440 and Choice RT, such that reaction time decreased across training (Fig. 4; Cue RT: $t(26.31) =$
441 $3.47, p = .002$; Choice RT: $t(27.02) = 2.51, p = .018$). There was no significant difference
442 between male and female rats (Cue RT: $t(37.89) = 0.62, p = .538$; Choice RT: $t(28.78) = -0.36, p$
443 $= .720$). For controls, Cue RT and Choice RT were not significantly different between rats
444 learning the 1D tasks and the 2D tasks (Cue RT: $t(26.96) = 2.09, p = .045$; Choice RT: $t(27.00) =$
445 $0.26, p = .796$). For rats with PL lesions, Cue RT was significantly larger than the controls for
446 rats learning the 1D tasks (Fig. 4A-B; $t(27.02) = 3.92, p < .001$; Fig. 3C), but not the 2D tasks
447 ($t(26.97) = 1.25, p = .223$). However, there were no significant group differences in Choice RT
448 (Figs. 4C-D; 1D tasks: $t(27.04) = 1.55, p = .133$; 2D tasks: $t(26.89) = 0.99, p = .329$). Together,

449 these results suggest that the rats with PL lesions learning the 1D tasks required more time to
450 evaluate each stimulus. However, there were no significant differences in the amount of time to
451 execute a category decision. These results are task-specific, which suggests that this impairment
452 is a consequence of the 1D tasks having both relevant and irrelevant stimulus information.

453

454 *PL lesions impaired choice confidence for rats learning 1D tasks*

455 We then examined the effect of PL lesions on touch separation, a measure of choice confidence
456 during the Cue phase (see Material and Methods). A linear mixed effects model (fixed effects:
457 group, training session, a quadratic function across sessions; random effects: intercept, slope, the
458 quadratic function; covariate: sex) examined touch separation for the third touch across training
459 sessions. First, there was a main effect of training session, such that touch separation increased
460 across sessions (Fig. 4; $t(27.02) = 4.71, p < .001$). There was no significant difference in touch
461 separation between male and female rats ($t(26.16) = -0.93, p = .360$) as well as controls learning
462 the 1D tasks and 2D tasks ($t(26.95) = 0.30, p = .840$). For rats with PL lesions, touch separation
463 was impaired for the rats learning the 1D tasks (Fig. 4E; $t(27.38) = 2.82, p = .009$), but not 2D
464 tasks (Fig. 4F; $t(26.96) = .53, p = .601$). These results support the role of PL in learning 1D tasks
465 and suggests that these rats were less confident in their category decisions.

466

467 *PL lesions impair category generalization for 1D tasks but not 2D task*

468 After category training, each rat was presented with five testing sessions to examine category
469 generalization. Testing distributions had identical category means as the training distributions but
470 had increased variance along the relevant dimension (or relevant axis for the 2D tasks) to sample
471 from novel portions of the stimulus space (Fig. 5A). We segregated the testing distributions into

472 three trial types: stimuli that overlapped with the training distributions (Trained), novel stimuli
473 farther from the category boundary (Distal), and novel stimuli closer to the category boundary
474 (Proximal).

475 Linear mixed effects models (fixed effects: group, trial type, a quadratic function; random
476 effects: intercept, slope, and the quadratic function; covariate: sex) examined accuracy, Cue RT,
477 Choice RT, and touch separation during testing sessions. Generally, performance was poorer for
478 Proximal stimuli compared to Trained stimuli, suggesting that the rats perceived stimuli closer to
479 the category boundary as more difficult (Broschard et al., 2019). Specifically, accuracy and
480 touch separation for Proximal stimuli were significantly lower than Trained stimuli, and Choice
481 RT for Proximal stimuli was significantly larger than Trained stimuli (accuracy: $t(52) = 8.22, p <$
482 $.001$; touch separation: $t(52) = 2.49, p = .016$; Choice RT: $t(52) = 2.76, p = .008$). Cue RT did not
483 differ significantly between Proximal stimuli and Trained stimuli ($t(52) = 2.0, p = .057$).

484 Conversely, rats could easily generalize to the Distal stimuli, and there were no significant
485 differences between Distal stimuli and Trained stimuli (accuracy: $t(52) = 1.96, p = .055$; Cue RT:
486 $t(52) = 0.94, p = .353$; Choice RT: $t(52) = 0.85, p = .400$; touch separation: $t(52) = .89, p = .377$).

487 Finally, there were no significant differences in all dependent measures between controls that
488 learned the 1D tasks and 2D tasks (Figs. 5B-E; Accuracy: $t(26) = 0.77, p = .448$; Cue RT:
489 $t(31.08) = 0.73, p = .470$; Choice RT: $t(30.23) = 0.33, p = .747$; touch separation: $t(38.54) = 0.04,$
490 $p = .966$).

491 PL lesions impaired accuracy and touch separation for rats that learned the 1D tasks
492 (Figs. 5B,E; accuracy: $t(26) = 2.51, p = .019$; touch separation: $t(38.54) = 2.95, p = .039$), but not
493 the 2D tasks (accuracy: $t(26) = 0.43, p = .667$; touch separation: $t(38.54) = 0.41, p = .684$).

494 Furthermore, Cue RT was significantly larger for rats with PL lesions that learned the 1D tasks,

495 but not the 2D tasks (Fig. 5C; $t(31.08) = 2.61, p = .014$; $t(31.08) = 0.72, p = .480$, respectively).
496 PL lesions did not affect Choice RT (Fig. 5D; 1D tasks: $t(30.23) = 0.27, p = .787$; 2D tasks:
497 $t(30.23) = 0.97, p = .341$). There were no significant interactions between trial types (all p s >
498 $.05$). There also were no significant differences between male and female rats (all p > $.05$).
499 Together, these results are consistent with the results from training. PL lesions impaired category
500 generalization for rats that learned the 1D tasks, but not the rats that learned the 2D tasks. Rats
501 with PL lesions learning the 1D tasks had lower accuracy, required more time to categorize each
502 stimulus, and had less confidence in their category decisions.

503

504 *Simple Discrimination*

505 After category generalization, rats were trained to learn a control discrimination task. The trial
506 sequence was identical to category training, except only two objects were presented (instead of
507 categories of stimuli; Fig. 6A). This procedure was added to ensure the PL lesions did not cause
508 general deficits that were not specific to categorization (i.e., motivational, perceptual, motor,
509 etc.). Using a 2x2 between ANOVA, there were no significant differences in the number of
510 sessions to reach the learning criterion across groups (Fig. 6B; $F(3,25) = .37, p > .05$). These
511 results support the conclusion that the observed impairments were specific to categorization.

512

513 *SUSTAIN modeling: PL affects selective attention and category representations*

514 Using the neural network SUSTAIN, we created three manipulations that simulated potential
515 functions of the PL (Love & Gureckis, 2007). Two of these manipulations disrupted SUSTAIN's
516 feature tuning mechanism, which learns to shift attention to category-relevant dimensions. These
517 included 1) lesioning the feature tuning mechanism so that attention weights are static across

518 training and 2) shuffling the attention weights before each trial so that attention was not
519 consistently directed towards category-relevant dimensions. The third manipulation tested the
520 prediction that PL lesions limited the ability to recruit new clusters; this was modeled by
521 increasing a cluster recruitment threshold parameter. Five models were created using
522 combinations of these manipulations (Fig. 7B & 7D). Each model was fit to the averaged group
523 data (Fig. 7B & 7D). These models were compared to a control model that assumed the lesions
524 had no effect on learning. The rats' behavior was best explained when we shuffled the attention
525 weights before each trial and increased the cluster recruitment threshold for the lesion groups
526 (Fig. 7C; Model 5). These results suggest that the PL is important for maintaining attention to
527 category-relevant dimensions as well as building category representations. All models produced
528 a better fit than the control model that assumed the lesions had no effect on learning.

529 We then examined the best fitting model in Figure 7D (Model 5) to ascertain how the
530 lesions affected the cluster representations. Figure 7E shows that, for the controls, SUSTAIN
531 recruited two clusters (one per category) to learn the 1D tasks, but multiple clusters (~3-5 per
532 category) to learn the 2D tasks (Broschard et al., 2020). These results suggest that 1D categories
533 are normally represented by single prototypes, whereas 2D categories are normally represented
534 by multiple exemplars (Posner & Keele, 1968; Nosofsky, 1986, respectively). Rats with PL
535 lesions recruited fewer clusters compared to controls to learn the 2D tasks, a direct consequence
536 of increasing the cluster recruitment threshold. These results imply that rats with PL lesions
537 learning the 2D tasks may have had sparser category representations compared to controls, even
538 if performance was intact across training (Figure 7E).

539 We then examined the feature tuning mechanism of the best-fitting model to characterize
540 how the PL lesions affected selective attention. Figure 7F demonstrates that 1D tasks were

541 learned by incrementally shifting attention towards the category-relevant dimension (Broschard
542 et al., 2020). Specifically, the attention weight of the category-relevant dimension increased
543 across training trials, whereas the attention weight to the category-irrelevant dimension
544 decreased across training trials. Importantly, this differentiation was much slower and reached
545 lower levels for rats with PL lesions (Fig. 7F). This finding verifies that shuffling the attention
546 weights across trials reduced selective attention by impairing the model's ability to maintain
547 attention to the relevant dimension. Conversely, the 2D tasks were learned by dividing attention
548 between stimulus dimensions (Fig. 7F; Broschard et al., 2020). The attention weights for both
549 dimensions were equivalent across training, a pattern that was consistent for both controls and
550 rats with PL lesions.

551

552 *PL lesions impair perceptual recency effects*

553 SUSTAIN was best fit to the averaged group data when it was assumed that the PL lesions
554 reduced the ability to update category representations. Here, we tested this prediction further by
555 examining category learning on a trial-by-trial basis. We predicted that if the PL is critical for
556 updating representations, then the PL lesions should also impair perceptual recency effects,
557 where the learner biases category decisions according to recent training experiences. To test this,
558 we binned the accuracy of training trials according to the perceived similarity between the
559 current exemplar and the most recent exemplar (see Materials and Methods). Then, we
560 subtracted the binned accuracies from iterations where trial order was randomized. Positive
561 recency scores indicate that accuracy was facilitated because of trial order, negative scores
562 indicate that accuracy was impaired because of trial order, and 0 indicates that trial order had no
563 effect on category accuracy.

564 For controls, trial order affected category learning and was modulated by stimulus
565 similarity (Fig. 8). One-sample *t*-tests were used to assess whether the perceptual recency scores
566 were significantly different from 0. For controls learning the 1D and 2D tasks, scores were
567 significantly larger than 0 if the current stimulus was perceptually similar to the previous
568 stimulus (i.e., above the median similarity; 1D tasks: $t(7) = 3.16, p = .016$; 2D tasks: $t(7) = 2.86,$
569 $p = .024$). Conversely, scores were significantly smaller than 0 if the current stimulus was
570 perceptually dissimilar from the previous stimulus (i.e., below the median similarity; 1D tasks:
571 $t(7) = 2.97, p = .021$; 2D tasks: $t(7) = 3.01, p = .020$). These results indicate that accuracy was
572 facilitated if the current stimulus was perceptually similar to the most recent exemplar, but
573 accuracy was impaired if the current stimulus was perceptually dissimilar from the most recent
574 exemplar. For rats with PL lesions, none of the perceptual recency scores were significantly
575 different from 0, indicating that trial order did not affect accuracy (Fig. 8; all $p > .05$). Together,
576 these results indicate that rats normally bias their decisions according to recent training
577 experiences, which implies that they regularly update category representations. This process is
578 effectively absent in rats with PL lesions. This finding supports the SUSTAIN modeling and
579 indicates that the PL is critical for updating category representations.

580

581 **Discussion**

582 Rats were trained to categorize stimuli containing black and white gratings according to one
583 stimulus dimension (1D tasks) or two dimensions (2D tasks). Lesions of the PL impaired
584 learning and generalization in rats trained on the 1D tasks. Without the PL, rats learning the 1D
585 tasks had lower accuracy, a larger number of correction trials, and more perseverative errors
586 compared to controls (Fig. 3); they also needed more time to categorize each stimulus (i.e., Cue

587 RT) and showed impaired choice confidence (i.e., touch separation; Fig. 4). The PL lesions did
588 not affect performance on the 2D tasks or the simple discrimination task; therefore, impairments
589 were specific to the 1D tasks. 1D and 2D tasks only differed in a simple rotation of the category
590 distributions. This rotation did not change any physical properties of the categories (Ashby,
591 Smith, & Rosedahl, 2019; e.g., discriminability, average category distance, etc.), but it did affect
592 how the tasks were learned by changing the number of category-relevant dimensions.

593 COVIS posits that humans have a PFC-mediated declarative system that learns new
594 categories by testing rules (Ashby et al., 1998; Ashby & Maddox, 2011). This system is biased
595 towards simple rules; therefore, COVIS predicts that the PFC is critical for learning tasks that
596 can be solved by unidimensional strategies (i.e., 1D tasks, but not 2D tasks). Using this logic, we
597 could conclude that rats also have a PFC-mediated declarative system that is important for
598 learning 1D tasks. However, there is little evidence that rats consistently apply category rules in
599 the manner that humans do (Broschard et al., 2019). Rule-based learning in humans is best
600 characterized by a step-wise learning curve, where accuracy improves rapidly in a non-linear
601 way (Ashby & Maddox, 2011). Presumably, this jump in performance is a consequence of the
602 participant testing hypotheses about potential rules and selecting the correct rule. Category
603 learning in rats is generally linear and incremental, even for the 1D tasks, suggesting that rats are
604 not testing hypotheses in the same way.

605 Instead, we propose that rodent PL mediates lower-level mechanisms that make up the
606 building blocks of the primate declarative system. Specifically, the rodent PL biases attention to
607 relevant stimulus information, a mechanism important for learning 1D tasks, but not for learning
608 2D tasks. This interpretation is supported by SUSTAIN. The neural network model best fit the
609 PL lesion data when we shuffled the attention weights before each decision, suggesting that the

610 PL normally maintains attention to relevant stimulus information (Fig. 7). Shuffling the attention
611 weights did not affect performance on the 2D tasks since attention was allocated to both
612 dimensions equally. This interpretation converges with multiple studies implicating the PL in
613 selective attention by orienting attention to cues that predict reward (Sharpe & Killcross, 2015,
614 2018; Tait et al., 2014).

615 Selective attention is foundational to categorization (Nosofsky, 1986). At its core,
616 category learning involves discriminating between relevant and irrelevant stimulus information.
617 To illustrate this point, Rehder & Hoffman (2005) tracked eye movements while participants
618 learned to categorize stimuli made from three binary dimensions; depending on the task, the
619 number of category-relevant dimension(s) differed (Shepard, Hovland, & Jenkins, 1961). Eye
620 fixations (and presumably attention) were initially distributed across all stimulus dimensions, but
621 then became restricted to only the relevant dimensions (Rehder & Hoffman, 2005). Our results
622 suggest that maintaining attention to a subset of stimulus dimensions is mediated by the PL, a
623 function that becomes more critical as the number of relevant dimensions decreases. This
624 interpretation also matches the results of Mack and colleagues (2020), who found that BOLD
625 activity in the ventromedial PFC (vmPFC) tracked the number of relevant stimulus dimensions.
626 They argued that the that vmPFC was critical for filtering out irrelevant stimulus information.

627 Future experiments can investigate whether other prefrontal subregions are also necessary
628 for learning 1D tasks. A potential target is the anterior cingulate cortex (ACC), which has also
629 been implicated in selective attention in rats (Kim, Wasserman, Castro, & Freeman, 2016).
630 COVIS posits that the ACC participates in the declarative system by switching attention to
631 alternative category rules (Ashby et al., 1998). This can be tested directly by inactivating the
632 rodent ACC before category training. One interesting prediction would be that the PL and ACC

633 serve similar but dissociable functions in selective attention. For example, whereas our results
634 suggest that the PL is critical for maintaining attention to relevant dimensions, the ACC may be
635 critical for identifying dimensions that are category-relevant vs. irrelevant. In this example, the
636 ACC would be critical for learning how to orient attention, and the PL would be critical in
637 applying those learned attention weights.

638 In addition to selective attention, the results from the SUSTAIN modeling suggest that
639 the PL is also important for creating new category representations (i.e., clusters). SUSTAIN
640 recruits new clusters in response to ‘surprising’ stimuli, where the model is confident in an
641 ultimately incorrect decision (Love et al., 2004). In the current experiment, SUSTAIN best fit the
642 learning data when it was assumed that the rats with PL lesions had a higher threshold to recruit
643 new clusters (Fig. 7). Consequently, without the PL, the category representations were much
644 sparser. This was especially critical for rats learning the 2D tasks, where normally multiple
645 clusters are recruited for each category. The role of the PL in updating category representations
646 was also examined by analyzing category learning on a trial-by-trial basis (Fig. 8). We found
647 that, for controls, category decisions were directly influenced by recent exemplars. Accuracy was
648 facilitated if the current stimulus was perceptually similar to the previous exemplar, whereas
649 accuracy was impaired if the current stimulus was dissimilar to the previous exemplar,
650 suggesting that rats update category decisions regularly and bias their decisions according to
651 recent information. Importantly, rats with PL lesions showed no effects of trial order. Without
652 the PL, rats may be less sensitive to local changes within the category, which could lead to
653 perseveration in the event of a task switch.

654 We predict that the role of the PL in updating representations is related to the literature
655 that credits the PFC in the development and maintenance of schemas, which are hierarchical

656 representations of information that help organize memories (Koscik & Tranel, 2012). Schemas
657 extrapolate common elements from distinct episodes (Morton, Sherrill, & Preston, 2017;
658 Pudhiyidath, Roome, Coughlin, Nguyen, & Preston, 2019) and rely on an interaction between the
659 PFC and hippocampus (Zeithamova, Dominick, & Preston, 2012; Schlichting & Preston, 2016).
660 We predict that the PL uses these mechanisms in our categorization tasks to update and elaborate
661 category representations. Indeed, a growing literature suggests that the hippocampus stores
662 category representations that are similar to the clusters described by SUSTAIN (Theves,
663 Fernandez, & Doeller, 2020; Mack, Love, & Preston, 2016; Mack, Love, & Preston, 2018). For
664 example, Mok & Love, 2020 was able to fit a clustering model to the neural activity of place
665 cells and grids cells as a rat navigated an environment. This implies that updating and building
666 category representations involves a close interaction between the PL and hippocampus. Future
667 experiments can examine this interaction directly.

668 Finally, it is important to note that although the PL facilitates category learning, it may
669 not be necessary for categorization to occur. Indeed, accuracy impairments in the 1D tasks
670 largely occurred during the initial training sessions, and rats with PL lesions were able to learn
671 the 1D tasks after extensive training. This implies that other neural regions were able to
672 compensate. COVIS predicts that a second learning system, the non-declarative system, takes
673 over when the PFC-mediated declarative system cannot successfully find a category rule (Ashby
674 et al., 1998; Ashby & Maddox, 2011). Importantly, key features of this non-declarative system
675 were present in rats with PL lesions. For instance, the non-declarative system does not employ
676 executive functions like selective attention. Additionally, learning in the non-declarative system
677 is thought to be more static and habitual, relying on repetition and consistent feedback. We
678 suspect that in the absence of the PL, a learning system synonymous to the non-declarative

679 system of COVIS compensated. We hypothesize that the dorsolateral striatum (the tail of the
680 caudate nucleus in primates) supports categorization in the absence of the PL, as this region is
681 important for supporting habitual behaviors in rats (Balleine, Delgado & Hikosaka, 2007).

682 To conclude, a general function of the PFC is to guide behaviors in an adaptive way
683 (Miller & Cohen, 2001). In the context of category learning, we conclude that the rodent PL
684 accomplishes this function through two mechanisms. First, the PL maintains attention to relevant
685 stimulus information (i.e., selective attention); this prevents the incorporation of irrelevant
686 information into category decisions. Second, the PL regularly updates category representations
687 and biases decisions according to recent information; this allows for dense, flexible
688 representations and primes the organism for changes in the category structure. Together, these
689 mechanisms allow for category representations that are both flexible and adaptive.

690

691

692

693

694

695

696

697

698

699

700

701

702 **References**

- 703 Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on*
704 *Automatic Control*, 19(6), 716–723. doi: 10.1109/tac.1974.1100705
- 705 Ashby, F. G., Alfonso-Reese, L. A., Turken, A. U., & Waldron, E. M. (1998). A
706 neuropsychological theory of multiple systems in category learning. *Psychological*
707 *Review*, 105(3), 442–481. <https://doi.org/10.1037//0033-295X.105.3.442>
- 708 Ashby, F. G., & Maddox, W. T. (2011). Human category learning 2.0: Human category learning
709 2.0. *Annals of the New York Academy of Sciences*, 1224(1), 147–161.
710 <https://doi.org/10.1111/j.1749-6632.2010.05874.x>
- 711 Ashby, F. G., Smith, J. D., & Rosedahl, L. A. (2019). Dissociations between rule-based and
712 information-integration categorization are not caused by differences in task difficulty.
713 *Memory & Cognition*, 48(4), 541–552. <https://doi.org/10.3758/s13421-019-00988-4>
- 714 Balleine, B. W., Delgado, M. R., & Hikosaka, O. (2007). The Role of the Dorsal Striatum in
715 Reward and Decision-Making. *Journal of Neuroscience*, 27(31), 8161–8165. doi:
716 10.1523/jneurosci.1554-07.2007
- 717 Broschard, M. B., Kim, J., Love, B. C., & Freeman, J. H. (2020). Category learning in rodents
718 using touchscreen-based tasks. *Genes, Brain and Behavior*, xxx. doi: 10.1111/gbb.12665
- 719 Broschard, M. B., Kim, J., Love, B. C., Wasserman, E. A., & Freeman, J. H. (2019). Selective
720 attention in rat visual category learning. *Learning & Memory*, 26(3), 84–92.
721 <https://doi.org/10.1101/lm.048942.118>
- 722 Bruin J. P., Sánchez-Santed F, Heinsbroek R. P., Donker A, & Postmes P. (1994). A behavioural
723 analysis of rats with damage to the medial prefrontal cortex using the morris water maze:

724 evidence for behavioural flexibility, but not for impaired spatial navigation. *Brain*
725 *Research*, 652(1), 323–333. doi:10.1016/0006-8993(94)90243-7

726 Bowman, C.R., Zeithamova, D. (2018). Abstract memory representations in the ventromedial
727 prefrontal cortex and hippocampus support concept generalization. *Journal of*
728 *Neuroscience*, 38(10), 2605-2614.

729 Crawley MJ. 2007. The R book. John Wiley & Sons Ltd., Chichester.

730 Crijns, E., & Op de Beeck, H. (2019). The Visual Acuity of Rats in Touchscreen Setups. *Vision*,
731 4(1), 4. <https://doi.org/10.3390/vision4010004>

732 Freedman, D. J. (2001). Categorical Representation of Visual Stimuli in the Primate Prefrontal
733 Cortex. *Science*, 291(5502), 312–316. <https://doi.org/10.1126/science.291.5502.312>

734 Heilbronner, S. R., Rodriguez-Romaguera, J., Quirk, G. J., Groenewegen, H. J., & Haber, S. N.
735 (2016). Circuit-Based Corticostriatal Homologies Between Rat and Primate. *Biological*
736 *Psychiatry*, 80(7), 509–521. doi: 10.1016/j.biopsych.2016.05.012

737 Horst N & Laubach M. (2009). The role of rat dorsomedial prefrontal cortex in spatial working
738 memory. *Neuroscience*, 164(1), 444–456. doi:10.1016/j.neuroscience.2009.08.004

739 Jones, M., Love, B. C., & Maddox, W. T. (2006). Recency effects as a window to generalization:
740 Separating decisional and perceptual sequential effects in category learning. *Journal of*
741 *Experimental Psychology: Learning, Memory, and Cognition*, 32(2), 316–332. doi:
742 10.1037/0278-7393.32.3.316

743 Kim, J., Castro, L., Wasserman, E. A., & Freeman, J. H. (2018). Dorsal hippocampus is
744 necessary for visual categorization in rats. *Hippocampus*, 28(6), 392–405.
745 <https://doi.org/10.1002/hipo.22839>

746 Kim, J., Wasserman, E. A., Castro, L., & Freeman, J. H. (2016). Anterior cingulate cortex
747 inactivation impairs rodent visual selective attention and prospective memory. *Behavioral*
748 *Neuroscience*, *130*(1), 75–90. <https://doi.org/10.1037/bne0000117>

749 Kosciak, T. R., & Tranel, D. (2012). The Human Ventromedial Prefrontal Cortex Is Critical for
750 Transitive Inference. *Journal of Cognitive Neuroscience*, *24*(5), 1191–1204. doi:
751 10.1162/jocn_a_00203

752 Kumaran, D., Summerfield, J. J., Hassabis, D., & Maguire, E. A. (2009). Tracking the
753 Emergence of Conceptual Knowledge during Human Decision Making. *Neuron*, *63*(6),
754 889–901. doi: 10.1016/j.neuron.2009.07.030

755 Laubach, M., Amarante, L. M., Swanson, K., & White, S. R. (2018). What, if anything, is rodent
756 prefrontal cortex? *Eneuro*, *5*(5). doi:10.1523/eneuro.0315-18.2018

757 Love, B. C., & Gureckis, T. M. (2007). Models in search of a brain. *Cognitive, Affective, &*
758 *Behavioral Neuroscience*, *7*(2), 90–108. doi: 10.3758/cabn.7.2.90

759 Love, B. C., Medin, D. L., & Gureckis, T. M. (2004). SUSTAIN: A Network Model of Category
760 Learning. *Psychological Review*, *111*(2), 309–332. [https://doi.org/10.1037/0033-](https://doi.org/10.1037/0033-295X.111.2.309)
761 [295X.111.2.309](https://doi.org/10.1037/0033-295X.111.2.309)

762 Mack, M. L., Love, B. C., & Preston, A. R. (2016). Dynamic updating of hippocampal object
763 representations reflects new conceptual knowledge. <https://doi.org/10.1101/071118>

764 Mack, M. L., Love, B. C., & Preston, A. R. (2018). Building concepts one episode at a time: The
765 hippocampus and concept formation. *Neuroscience Letters*, *680*, 31-38.
766 doi:10.1016/j.neulet.2017.07.061

767 Mack, M. L., Preston, A. R., & Love, B. C. (2020). Ventromedial prefrontal cortex compression
768 during concept learning. *Nature Communications*, *11*(1), 1-11. doi: 10.1038/s41467-019-

769 13930-8

770 Maddox, W. T., Ashby, F. G., & Bohil, C. J. (2003). Delayed feedback effects on rule-based and
771 information-integration category learning. *Journal of Experimental Psychology:*
772 *Learning, Memory, and Cognition*, 29(4), 650–662. doi: 10.1037/0278-7393.29.4.650

773 Marquis J, Killcross S, Haddon J. E. (2007). Inactivation of the prelimbic, but not infralimbic,
774 prefrontal cortex impairs the contextual control of response conflict in rats. *European*
775 *Journal of Neuroscience*, 25(1), 559–566. doi:10.1111/j .1460-9568.2006.05295.x

776 Miller, E. K., & Cohen, J. D. (2001). An Integrative Theory of Prefrontal Cortex
777 Function. *Annual Review of Neuroscience*, 24(1), 167–202. doi:
778 10.1146/annurev.neuro.24.1.167

779 Morton, N. W., Sherrill, K. R., & Preston, A. R. (2017). Memory integration constructs maps of
780 space, time, and concepts. *Current Opinion in Behavioral Sciences*, 17, 161–168. doi:
781 10.1016/j.cobeha.2017.08.007

782 Nomura, E., Maddox, W., Filoteo, J., Ing, A., Gitelman, D., Parrish, T., Mesulam, M.-M., &
783 Reber, P. (2006). Neural Correlates of Rule-Based and Information-Integration Visual
784 Category Learning. *Cerebral Cortex*, 17(1), 37–43. <https://doi.org/10.1093/cercor/bhj122>

785 Nosofsky, R. M. (1986). Attention, similarity, and the identification-categorization relationship.
786 *Journal of Experimental Psychology: General*, 115(1), 39–57.
787 <https://doi.org/10.1037//0096-3445.115.1.39>

788 O’Donoghue, E. M., Broschard, M. B., & Wasserman, E. A. (2020). Pigeons exhibit flexibility
789 but not rule formation in dimensional learning, stimulus generalization, and task
790 switching. *Journal of Experimental Psychology: Animal Learning and Cognition*, 46(2),
791 107–123. doi: 10.1037/xan0000234

792 Ostlund S.B. (2005). Lesions of medial prefrontal cortex disrupt the acquisition but not the
793 expression of goal-directed learning. *Journal of Neuroscience*, 25(1), 7763–7770.
794 doi:10.1523/JNEUROSCI.1921-05.2005

795 Posner, M. I., & Keele, S. W. (1968). On the genesis of abstract ideas. *Journal of Experimental*
796 *Psychology*, 77(3), 353–363. doi: 10.1037/h0025953

797 Pudhiyidath, A., Roome, H. E., Coughlin, C., Nguyen, K. V., & Preston, A. R. (2019).
798 Developmental differences in temporal schema acquisition impact reasoning decisions.
799 *Cognitive Neuropsychology*, 1–21. doi: 10.1080/02643294.2019.1667316

800 Ragozzino M. E. (2007). The contribution of the medial prefrontal cortex, orbitofrontal cortex,
801 and dorsomedial striatum to behavioral flexibility. *Annals of the New York Academy of*
802 *Sciences*, 1121(1), 355–375. doi:10.1196/annals.1401.013

803 Rehder, B., & Hoffman, A. B. (2005). Eyetracking and selective attention in category learning.
804 *Cognitive Psychology*, 51(1), 1–41. <https://doi.org/10.1016/j.cogpsych.2004.11.001>

805 Reinert, S., Hübener, M., Bonhoeffer, T., & Goltstein, P. M. (2021). Mouse prefrontal cortex
806 represents learned rules for categorization. *Nature*. [https://doi.org/10.1038/s41586-021-](https://doi.org/10.1038/s41586-021-03452-z)
807 [03452-z](https://doi.org/10.1038/s41586-021-03452-z)

808 Rose, J. E., & Woolsey, C. N. (1948). Structure and relations of limbic cortex and anterior
809 thalamic nuclei in rabbit and cat. *The Journal of Comparative Neurology*, 89(3), 279–
810 347. doi: 10.1002/cne.900890307

811 Seamans J. K., Lapish C. C., & Durstewitz D. (2008). Comparing the prefrontal cortex of rats
812 and primates: insights from electrophysiology. *Neurotoxicity Research*, 14(1), 249–262.
813 doi:10.1007/BF03033814

814 Schlichting, M. L., & Preston, A. R. (2016). Hippocampal–medial prefrontal circuit supports
815 memory updating during learning and post-encoding rest. *Neurobiology of Learning and*
816 *Memory*, 134, 91–106. doi: 10.1016/j.nlm.2015.11.005

817 Sharpe, M. J., & Killcross, S. (2015). The prelimbic cortex directs attention toward predictive
818 cues during fear learning. *Learning & Memory*, 22(6), 289–293.
819 <https://doi.org/10.1101/lm.038273.115>

820 Sharpe, M. J., & Killcross, S. (2018). Modulation of attention and action in the medial prefrontal
821 cortex of rats. *Psychological Review*, 125(5), 822–843.
822 <https://doi.org/10.1037/rev0000118>

823 Shepard, R. N., Hovland, C. I., & Jenkins, H. M. (1961). Learning and memorization of
824 classifications. *Psychological Monographs: General and Applied*, 75(13), 1-42.
825 doi:10.1037/h0093825

826 Smith, J. D., Berg, M. E., Cook, R. G., Murphy, M. S., Crossley, M. J., Boomer, J., ... Grace, R.
827 C. (2012). Implicit and explicit categorization: A tale of four species. *Neuroscience &*
828 *Biobehavioral Reviews*, 36(10), 2355–2369. doi: 10.1016/j.neubiorev.2012.09.003

829 Tait D. S., Bowman E. M., Neuwirth L. S., & Brown V. J. (2018). Assessment of
830 intradimensional/extradimensional attentional set-shifting in rats. *Neuroscience*
831 *Biobehavioral Review*, 89(1), 72–84. doi:10.1016/j.neubiorev.2018.02.013

832 Tait D. S., Chase E., & Brown V. (2014). Attentional set-shifting in rodents: a review of
833 behavioural methods and pharmacological results. *Current Pharmaceutical Design*,
834 20(1), 5046–5059. doi:10.2174/1381612819666131216115802

835

836

837 Theves, S., Fernández, G., & Doeller, C. F. (2020). The Hippocampus Maps Concept Space, Not
838 Feature Space. *The Journal of Neuroscience*, 40(38), 7318–7325.
839 <https://doi.org/10.1523/jneurosci.0494-20.2020>

840 Uylings H. B., & Eden C.G. (1991). Chapter 3 Qualitative and quantitative comparison of the
841 prefrontal cortex in rat and in primates, including humans. *Progress in Brain Research*,
842 85(1), 31–62. doi:10.1016/s0079-6123(08)62675-8

843 van Aerde K. I., & Feldmeyer D. (2013). Morphological and physiological characterization of
844 pyramidal neuron subtypes in rat medial prefrontal cortex. *Cerebral Cortex*, 25(1), 788–
845 805. doi:10.1093/cercor/bht278

846 Vinken, K., Vermaercke, B., & Op de Beeck, H. P. (2014). Visual Categorization of Natural
847 Movies by Rats. *Journal of Neuroscience*, 34(32), 10645–10658.
848 <https://doi.org/10.1523/jneurosci.3663-13.2014>

849 Wit, S. D., Kosaki Y, Balleine B. W., & Dickinson A. (2006). Dorsomedial Prefrontal Cortex
850 Resolves Response Conflict in Rats. *Journal of Neuroscience*, 26(19), 5224-5229.
851 doi:10.1523/jneurosci.5175-05.2006

852 Yoon T, Okada J, Jung M. W., & Kim J.J. (2008). Prefrontal cortex and hippocampus subserve
853 different components of working memory in rats. *Learning & Memory*, 15(1), 97–105.
854 doi:10.1101/lm.850808

855 Zeithamova, D., Dominick, A. L., & Preston, A. R. (2012). Hippocampal and Ventral Medial
856 Prefrontal Activation during Retrieval-Mediated Learning Supports Novel Inference.
857 *Neuron*, 75(1), 168–179. doi: 10.1016/j.neuron.2012.05.010

858 Zoccolan, D., Oertelt, N., DiCarlo, J. J., & Cox, D. D. (2009). A rodent model for the study of
859 invariant visual object recognition. *Proceedings of the National Academy of Sciences*,
860 106(21), 8748–8753. <https://doi.org/10.1073/pnas.0811583106>

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

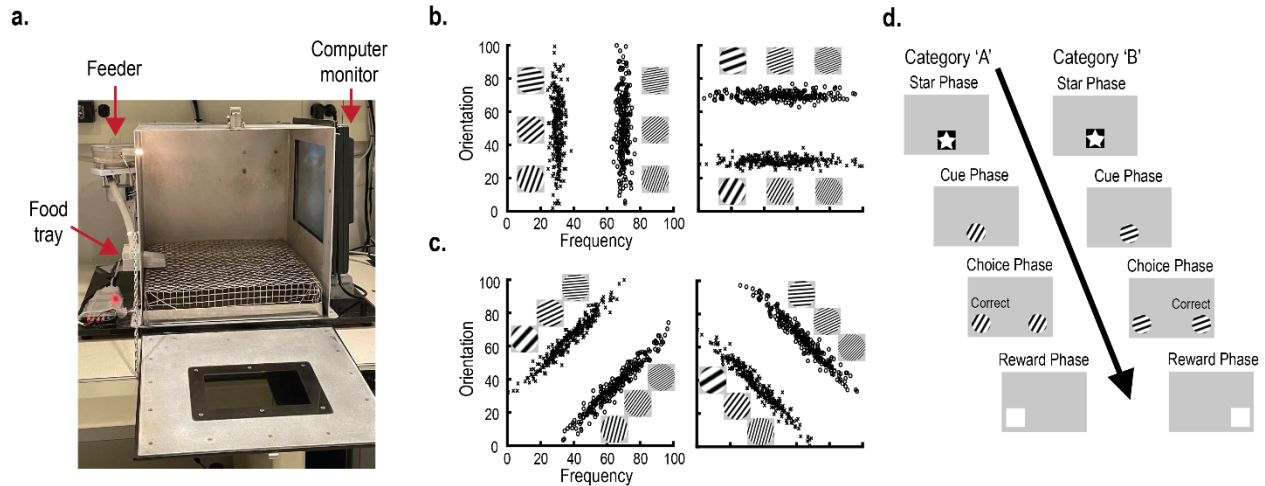
876

877

878

879

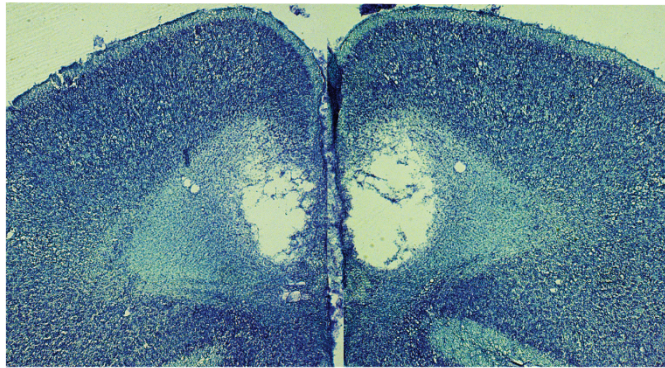
880



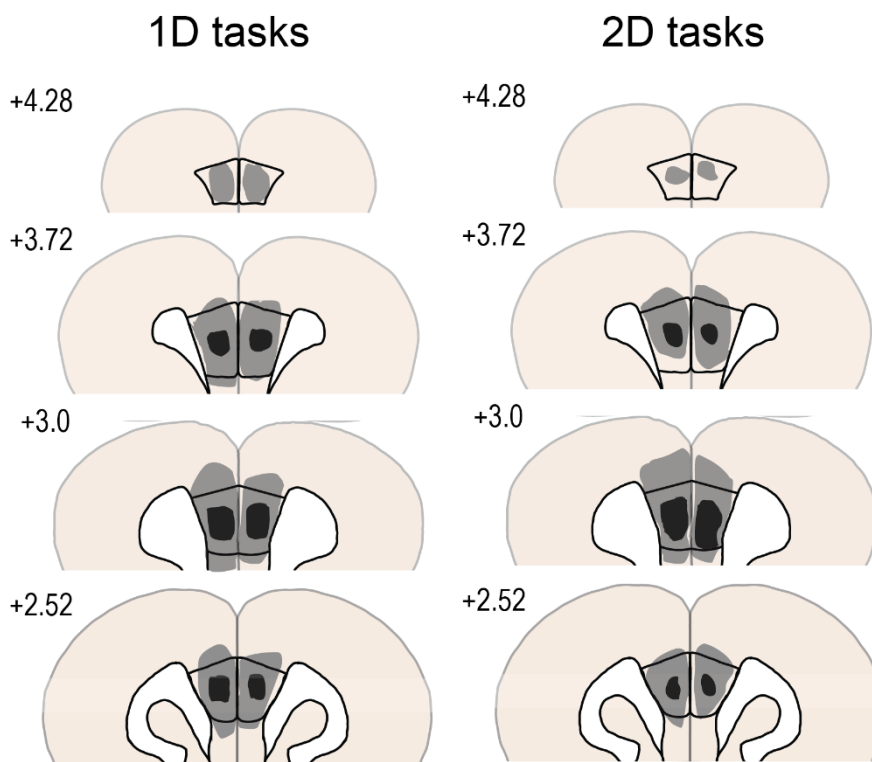
881 **Figure 1. A,** Behavioral testing was conducted in custom-built chambers. Each chamber
 882 contained a computer monitor and a touchscreen panel so that the rats could interact with the
 883 visual stimuli. A feeder delivered food pellets into a food tray to reinforce behavior. **B-C,** Rats
 884 were randomly assigned to learn one of four category tasks. For each task, category exemplars
 885 contained gratings that varied in their spatial frequency and orientation. Categories were created
 886 by placing normal distributions on this two-dimensional stimulus space. **B,** For the 1D tasks,
 887 category distributions were perpendicular to a stimulus axis. Consequently, one stimulus
 888 dimension was category-relevant (i.e., the dimension perpendicular to the distributions); the
 889 second dimension was category-irrelevant. We predicted that would rats use selective attention
 890 to learn 1D tasks by shifting attention towards the relevant dimension. **B,** For the 2D tasks,
 891 category distributions were not perpendicular to a stimulus axis. Therefore, both stimulus
 892 dimensions were category-relevant. **C,** The typical trial sequence for all training and testing
 893 sessions. Rats initiated each trial by touching the star stimulus at the center of the screen (Star
 894 phase). Then, an exemplar was randomly generated from the category distributions and placed at
 895 the center of the screen (Cue phase). The rat touched this exemplar three times, at which point
 896 copies of the exemplar were presented at the left and right sides of the screen (Choice phase).
 897 These copies acted as report keys. Members of category ‘A’ required a touch to the left report
 898 key, and members of category ‘B’ required a touch to the right report key. For correct responses,
 899 a white box appeared on the screen (Reward phase); one touch of the white box delivered a food
 900 reward. For incorrect responses, a correction trial was initiated, where the trial repeated from the
 901 Cue phase after a timeout.

902
 903
 904
 905
 906
 907
 908
 909
 910
 911
 912
 913
 914

a.

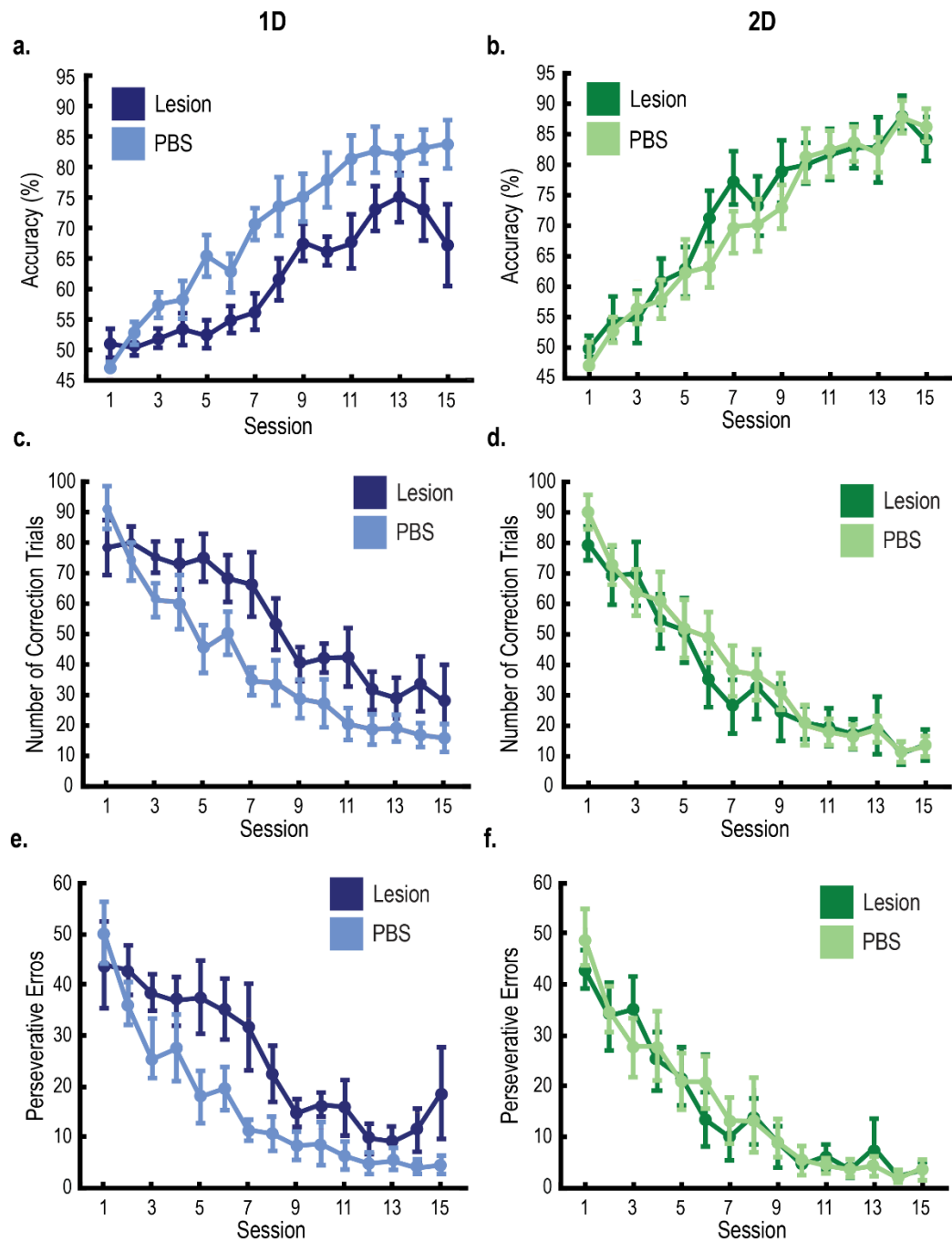


b.

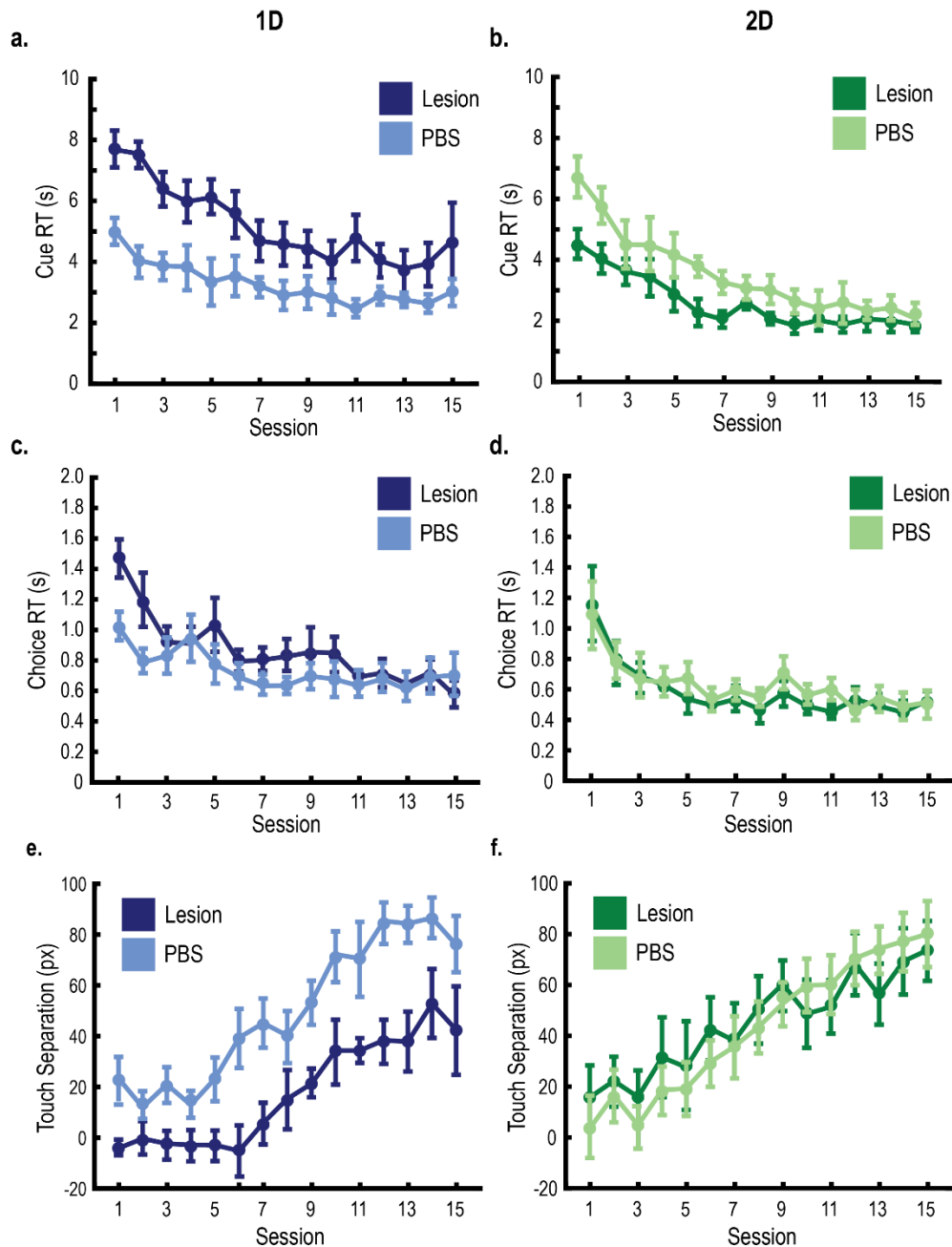


915
916
917
918
919
920
921
922
923
924
925
926

Figure 2. **A,** A representative example of the location and spread of the PL lesions. **B,** A comparison of lesion size and location for the smallest lesion (light gray) and the largest lesion (dark gray) for rats learning a 1D task (left) and rats learning a 2D task (right). All lesions were centered in the PL and were contained within bregma +4.3 and +2.2. Lesions rarely extended into cingulate cortex and infralimbic cortex.



927
 928 **Figure 3.** Excitotoxic lesions of the PL impaired learning 1D tasks, but not 2D tasks. **A-B**, Mean
 929 session accuracy of rats learning 1D tasks (A) and 2D tasks (B) (n = 8 per group). Compared to
 930 controls, rats with PL lesions had impaired accuracy for 1D tasks, but not for 2D tasks.
 931 Impairments were greatest at the beginning of category training. **C-D**, Mean number of
 932 correction trials from rats learning 1D tasks (C) and 2D tasks (D). Compared to controls, rats
 933 with PL lesions learning the 1D tasks, but not the 2D tasks required more correction trials. **E-F**,
 934 Mean number of perseverative errors for rats learning the 1D tasks (E) and 2D tasks (F).
 935 Compared to controls, rats with PL lesions learning the 1D tasks, but not the 2D tasks made
 936 more perseverative errors, where a choice was repeated after receiving negative feedback. All
 937 error bars indicate the *SEM*.



938

939

940

941

942

943

944

945

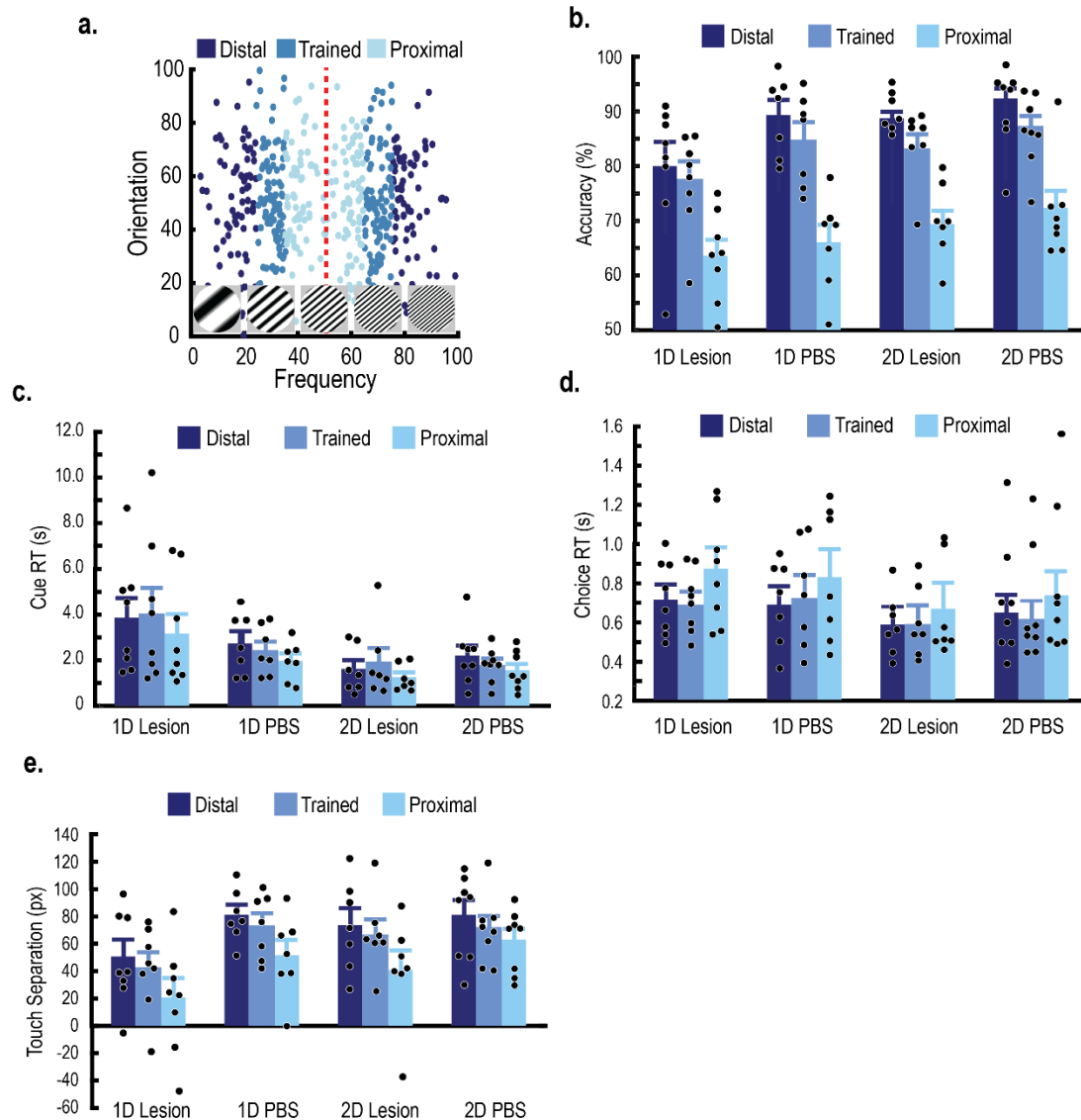
946

947

948

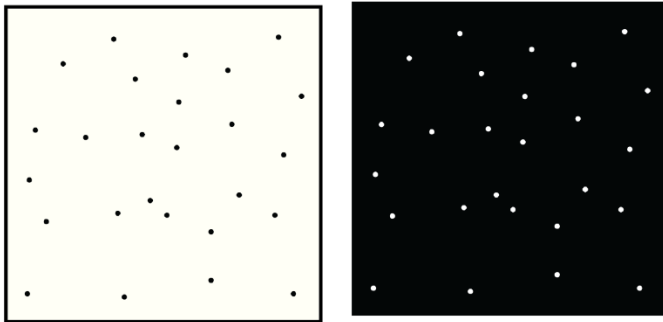
949

Figure 4. Excitotoxic lesions of the PL affected reaction time and choice anticipation during category learning. **A-B**, Mean time to observe and categorize each exemplar (Cue RT) for rats learning 1D tasks (A) and 2D tasks (B). Compared to controls, rats with PL lesions learning the 1D tasks, but not the 2D tasks exhibited a longer Cue RT. **C-D**, Mean time to execute a category decision (Choice RT) for rats learning the 1D tasks (C) and 2D tasks (D). Compared to controls, PL lesions did not affect Choice RT. **E-F**, Touch separation used the x-coordinate of the three touches during the Cue phase to estimate choice confidence. Positive touch separation indicates horizontal movement of the rat towards the correct side, whereas negative touch separation indicates horizontal movement towards the incorrect side. Compared to controls, rats with PL lesions learning the 1D tasks (A), but not the 2D tasks (B) exhibited lower touch separation across category learning. All error bars indicate the *SEM*.

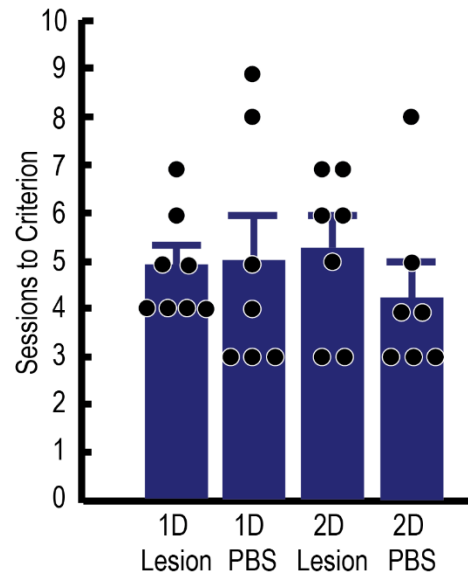


950
 951 **Figure 5.** The PL lesions impaired category generalization in rats trained on the 1D tasks, but not
 952 the 2D tasks. **A**, Each rat was given five testing sessions to examine category generalization.
 953 Testing distributions had the same category means as the training distributions, but the standard
 954 deviation along the relevant dimension was expanded to cover novel portions of the stimulus
 955 space. Each dot within the distributions represents a unique Gabor patch presented during
 956 testing. Testing distributions were split into three trial types: exemplars that overlapped with the
 957 training distributions (Trained), novel exemplars closer to the category boundary (Proximal), and
 958 novel exemplars farther from the category boundary (Distal). **B**, Mean accuracy across trial
 959 types. Generally, accuracy increased according to the distance from the category boundary. PL
 960 lesions impaired generalization in rats that learned the 1D tasks, but not rats that learned the 2D
 961 tasks. **C**, Mean Cue RT across trial types. Cue RT was larger for rats with PL lesions and had
 962 learned the 1D tasks than all other groups. There were no significant interactions across trial
 963 types. **D**, Mean Choice RT across trial types. Generally, Choice RT was larger for Proximal
 964 trials. The PL lesions did not affect Choice RT. **E**, Mean touch separation across trial types.
 965 Touch separation was reduced for rats with PL lesions that learned the 1D tasks. There were no
 966 significant interactions across trial types. All error bars indicate the *SEM*.

a.

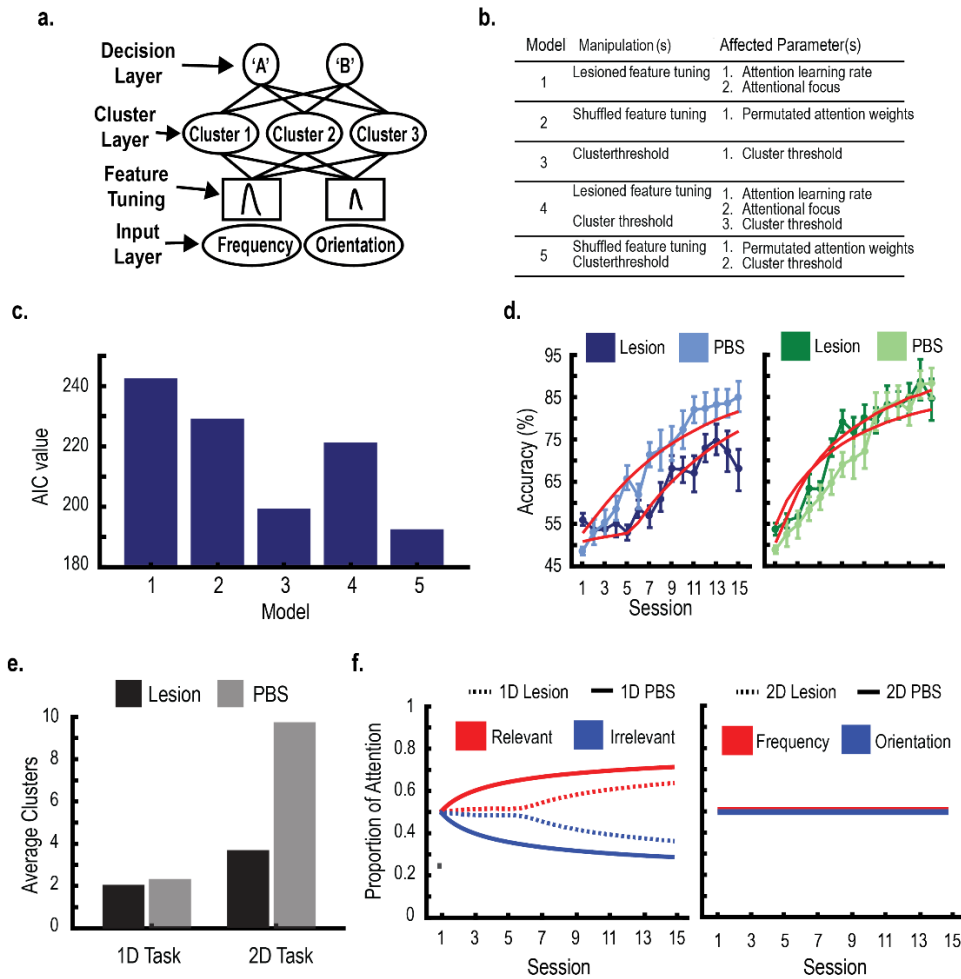


b.

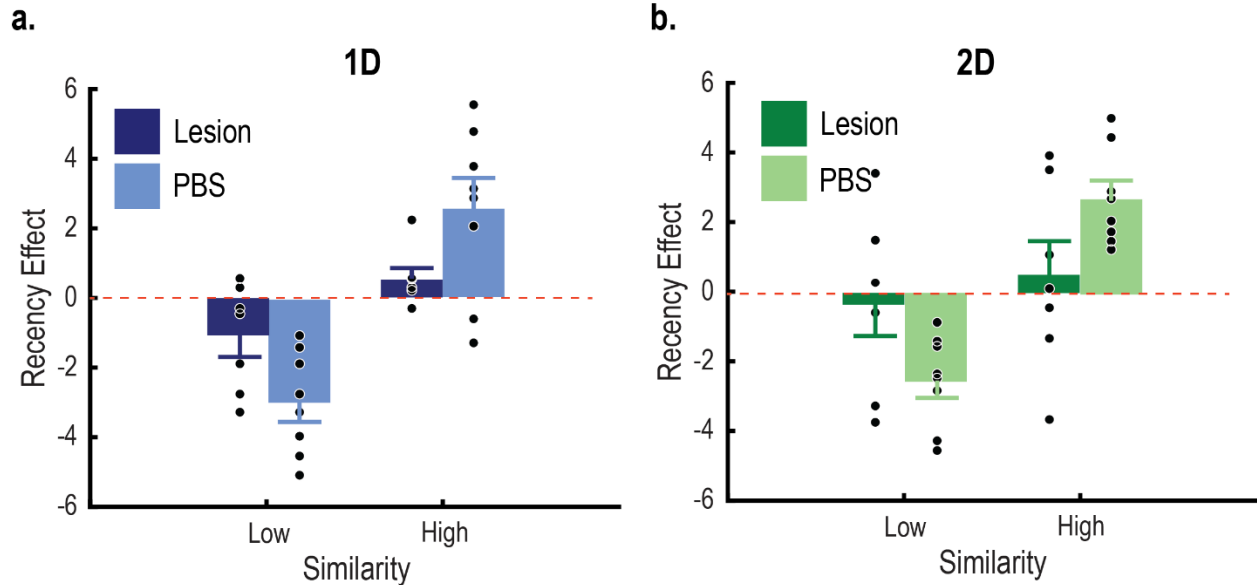


967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996

Figure 6. Rats were presented training sessions to learn to discriminate a dark box from a light box. All groups reached learning criterion (75% accuracy for both stimuli) in an equal number of training sessions. All error bars indicate the *SEM*.



997
 998 **Figure 7. A**, A diagram of the neural network model SUSTAIN, which contains three distinct
 999 layers: the input layer, cluster layer, and decision layer. SUSTAIN also contains a mechanism of
 1000 selective attention (i.e., the feature tuning mechanism) that weights stimulus information
 1001 according to category relevance. **B**, Descriptions of the five SUSTAIN models that were fit to
 1002 the learning data to test the effects of the PL lesions on category learning. These models were
 1003 compared to a control model which assumed the lesions had no effect on learning. **C**, The best
 1004 fitting model was determined by comparing the estimated AIC values. The model that best fit the
 1005 data (Model 5) assumed that the PL maintains attention to category-relevant information and
 1006 updates category representations. All models produced a better fit than the control model that
 1007 assumed the lesions had no effect on learning (not graphed: AIC = 278). **D**, SUSTAIN's
 1008 predictions using the best fitting model for rats learning the 1D (left) and 2D tasks (right). All
 1009 error bars indicate the *SEM*. **E**, Mean number of clusters recruited by SUSTAIN using the best
 1010 fitting model. Generally, SUSTAIN recruited two clusters (one per category) to learn the 1D
 1011 tasks and multiple clusters (3-4 per category) to learn the 2D tasks. For the rats with PL lesions,
 1012 the number of recruited clusters was reduced. **F**, The feature tuning mechanism of the best fitting
 1013 model. For rats learning the 1D tasks, the attention weight for the relevant dimension increased
 1014 across training, whereas the attention weight for the irrelevant dimension decreased across
 1015 training. This differentiation was impaired for rats with PL lesions. For rats learning the 2D
 1016 tasks, the attention weights were equivalent between dimensions and across training. This was
 1017 true for both control and lesioned rats.



1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025
 1026

Figure 8. Perceptual recency effects. Accuracy was binned according to the perceptual similarity between the current exemplar and the most recent exemplar. Then, these binned accuracies were subtracted from iterations where trial order was randomized. For controls learning both task types, accuracy was facilitated if the current stimulus had high perceptual similarity to the previous trial (i.e., a positive recency score). Accuracy was impaired if the current stimulus had low perceptual similarity to the previous trial (i.e., a negative recency score). These effects of trial order were absent in rats with PL lesions. This was true for rats learning the 1D (A) and 2D tasks (B). All error bars indicate the *SEM*.