

1 Behavioral flexibility is manipulatable and it improves flexibility
2 and problem solving in a new context.

3 Logan CJ^{1*} Lukas D^{1*} Blaisdell AP² Johnson-Ulrich Z³ MacPherson M³
4 Seitz B² Sevchik A⁴ McCune KB³

5 2022-08-15

6 Open...  access  code  peer review  data

7
8 **Affiliations:** 1) Max Planck Institute for Evolutionary Anthropology, Leipzig, Germany, 2) University of
9 California Los Angeles, USA, 3) University of California Santa Barbara, USA, 4) Arizona State University,
10 Tempe, AZ USA. *Corresponding author: corina_logan@eva.mpg.de

11
12 **This is the post-study manuscript of the preregistration that was pre-study peer reviewed and**
13 **received an In Principle Recommendation on 26 Mar 2019 by:**

14 Aurélie Coulon (2019) Can context changes improve behavioral flexibility? Towards a better un-
15 derstanding of species adaptability to environmental changes. *Peer Community in Ecology*, 100019.
16 [10.24072/pci.ecology.100019](https://doi.org/10.24072/pci.ecology.100019). Reviewers: Maxime Dahirel and Andrea Griffin

17 **Preregistration:** [html](#), [pdf](#), [rmd](#)

18 **Post-study manuscript** (submitted to PCI Ecology for post-study peer review on 3 Jan 2022): preprint
19 [pdf](#) at EcoEvoRxiv. Revised 15 Aug 2022: [pdf](#) at EcoEvoRxiv, [html](#), [rmd](#)

20 **ABSTRACT**

21 Behavioral flexibility, the ability to adapt behavior to new circumstances, is thought to play an important
22 role in a species' ability to successfully adapt to new environments and expand its geographic range. However,
23 flexibility is rarely directly tested in species in a way that would allow us to determine how flexibility works
24 to predict a species' ability to adapt their behavior to new environments. We use great-tailed grackles
25 (*Quiscalus mexicanus*; a bird species) as a model to investigate this question because they have recently
26 rapidly expanded their range into North America. We attempted to manipulate grackle flexibility using
27 colored tube reversal learning to determine whether flexibility is generalizable across contexts (multi-access
28 box), and what learning strategies grackles employ. We found that flexibility was manipulatable: birds in the
29 manipulated group took fewer trials to pass criterion with increasing reversal number, and they reversed a
30 color preference in fewer trials by the end of their serial reversals compared to control birds who had only one
31 reversal. Birds that passed their last reversal faster were also more flexible (faster to switch between loci) and
32 innovative (solved more loci) on a multi-access box. All grackles in the manipulated reversal learning group
33 used one learning strategy (epsilon-decreasing: long exploration period) in all reversals and did not use the
34 epsilon-first strategy: quickly shift their preference), and none used a particular exploration or exploitation
35 strategy earlier or later in their serial reversals. Understanding how flexibility causally relates to other traits

36 will allow researchers to develop robust theory about what flexibility is and when to invoke it as a primary
37 driver in a given context, such as a rapid geographic range expansion.

38 [Video summary](#)

39 INTRODUCTION

40 Behavioral flexibility, the ability to adapt behavior to new circumstances (see Mikhalevich et al., 2017 for
41 the theoretical background on this definition), is thought to play an important role in a species' ability
42 to successfully adapt to new environments and expand its geographic range (e.g., Lefebvre et al., 1997;
43 Sol et al., 2002, 2005, 2007; Sol & Lefebvre, 2000). The behavioral flexibility (hereafter referred to as
44 flexibility) of individuals is considered an important trait that facilitates the capacity for learning, which
45 is then associated with problem solving ability (applying what one has learned about the world to then
46 attempt to access a resource that is not readily accessible) (see review in Lea et al., 2020). It is hypothesized
47 that, through flexibility, individuals can increase the diversity of their behaviors either via asocial learning
48 (innovativeness) or social learning, leading to the establishment of the population in a new area (Wright et
49 al., 2010).

50 It is predicted that flexibility should positively relate with innovativeness, the ability to create a new behavior
51 or use an existing behavior in a new situation (Griffin & Guez, 2014). However, these predictions are based
52 on species-level data and proxies for flexibility and for innovation when examining such relationships (see
53 Logan et al., 2018). Flexibility is rarely directly tested in species that are rapidly expanding their geographic
54 ranges in a way that would allow us to determine how flexibility works and predict a species' ability to adapt
55 their behavior to new areas. Those investigations that examine the relationship between flexibility and
56 innovation [or problem solving - a type of experimental assay that does not necessarily require innovativeness
57 to solve, e.g., the ability to solve tasks using pre-trained behaviors; Griffin & Guez (2014)] in species that
58 are expanding their range show mixed results, with these variables correlating positively (e.g., grey squirrels:
59 Chow et al., 2016), negatively (e.g., Indian mynas: Griffin et al., 2013), or not at all (e.g., stick tool use and
60 string pulling in great-tailed grackles: Logan, 2016).

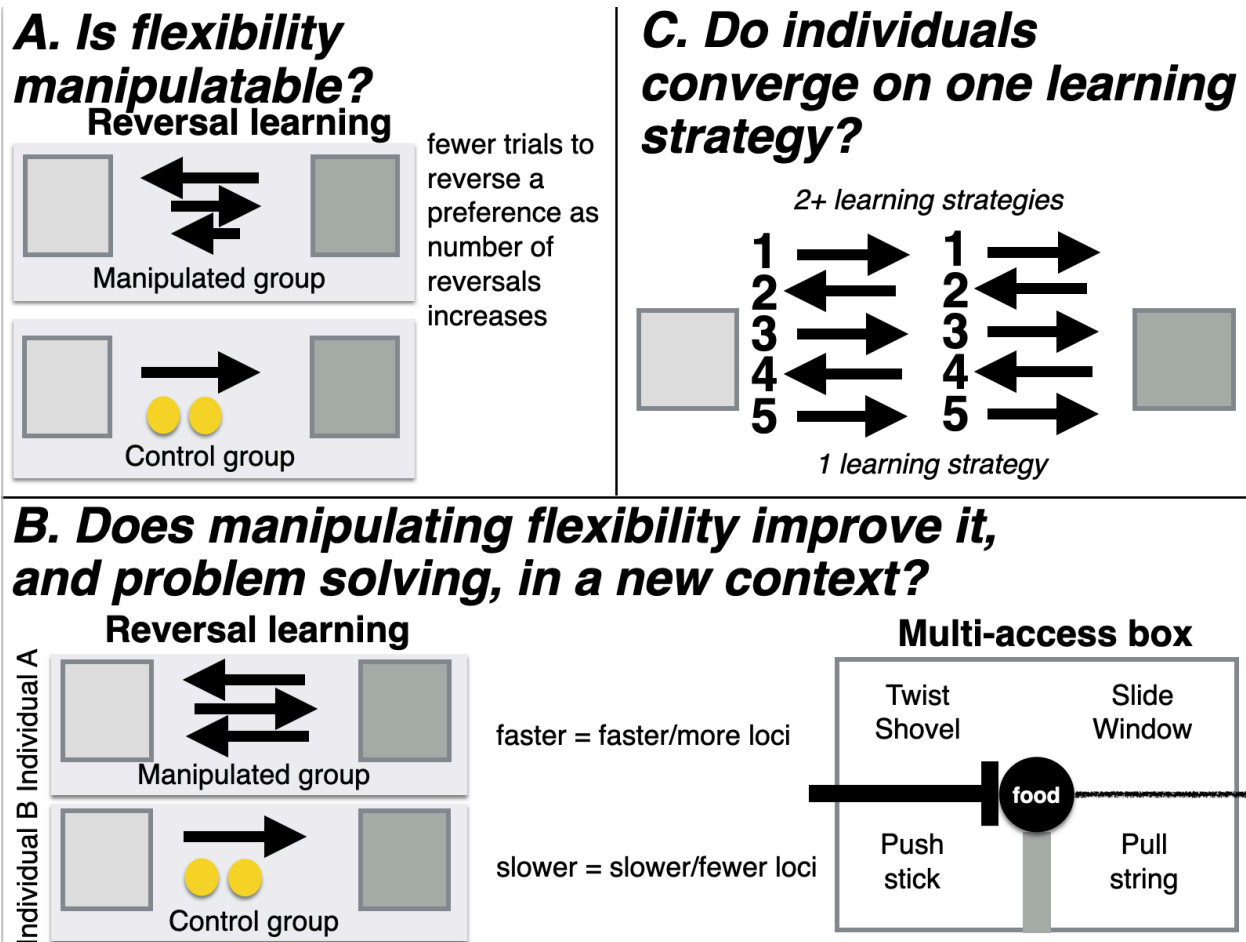
61 The first step to improving our understanding of whether and how flexibility relates to innovativeness, and
62 the focus of the current investigation, is to start with one population and perform a manipulative experiment
63 on one of the variables to determine whether there is an associated change in the other. Once this association
64 is known, future research can then investigate whether flexibility and innovativeness are involved in a range
65 expansion. Manipulative experiments go beyond correlations to infer a cause and effect relationship between
66 the manipulated variable and the variable(s) measured after the manipulation (Hernán & Robins, 2006;
67 McElreath, 2020). A manipulative experiment combined with the random assignment of subjects to a
68 condition (manipulated group or control group), eliminates many confounds associated with internal and
69 external variation (for example, season, motivation, sex, and so on). Such manipulative experiments in
70 behavioral ecology have primarily been conducted in laboratory settings because of the increased feasibility,
71 however such experiments are now also being conducted in wild settings (Aplin et al., 2015).

72 We focused our study on one population of great-tailed grackles (*Quiscalus mexicanus*, hereafter grackles),
73 a bird species that is flexible (Logan, 2016) and, while they are originally from Central America, they have
74 rapidly expanded their geographic range across the US since 1880 (Summers et al., 2022; Wehtje, 2003). We
75 attempted to manipulate grackle flexibility using serial reversals of a color preference to determine whether
76 their flexibility is generalizable across additional experimental contexts (touchscreen reversal learning and
77 multi-access box solution switching), whether improving flexibility also improves innovativeness (number of
78 loci solved on a multi-access box), and what learning strategies grackles employ (Figure 1).

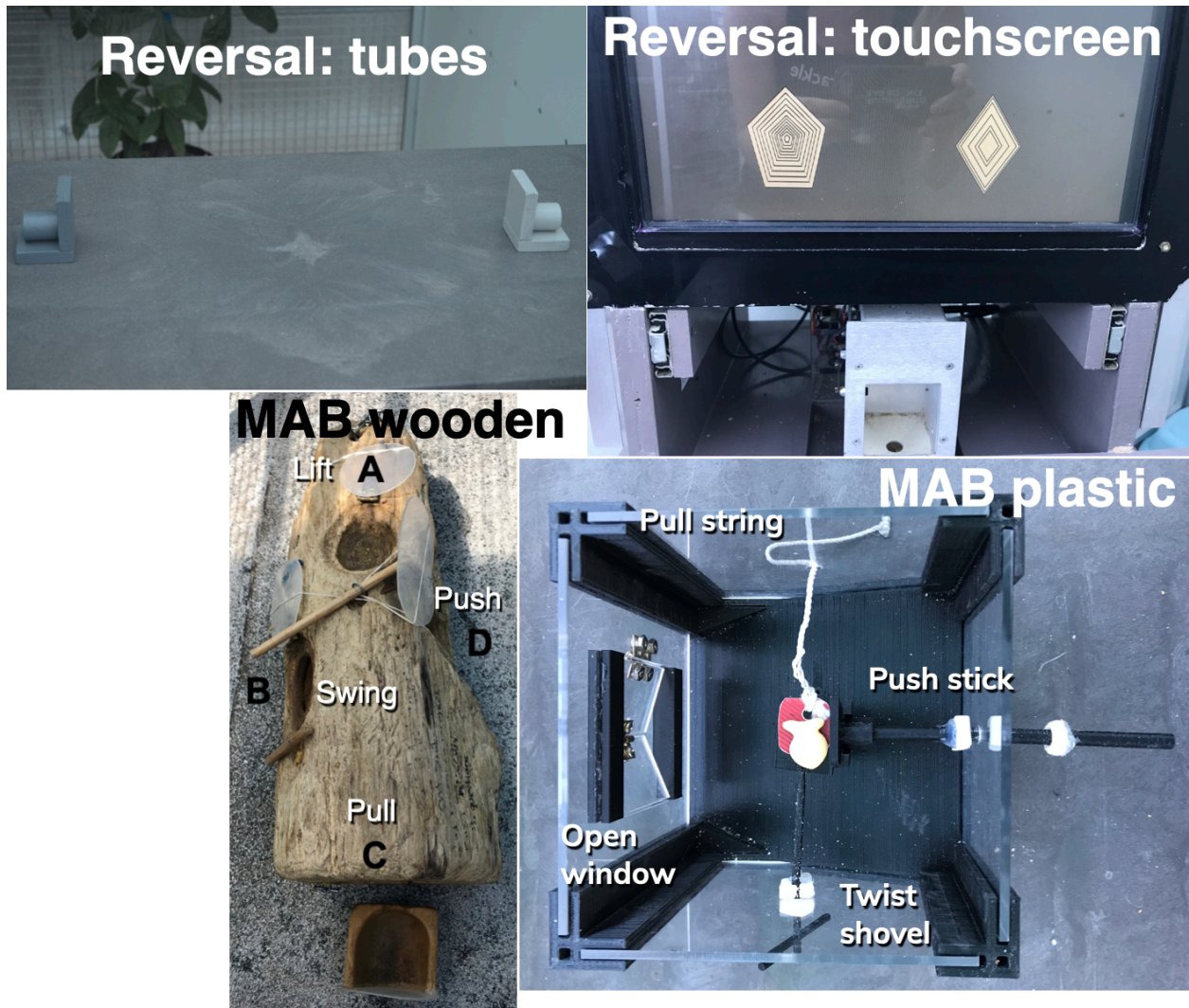
79 Reversal learning is a common way of measuring flexibility that has been used for many decades across many
80 species, therefore lending itself well to comparative analyses and generalizations (see review in Lea et al.,
81 2020). In this test, an individual learns to prefer the rewarded option, which differs from the non-rewarded
82 option in color, shape, space, or another obvious feature. Once this initial preference is formed, the previously
83 non-rewarded option becomes the rewarded option and vice versa, and the preference is reversed. Individuals

84 who are faster to reverse their preference are considered more flexible - better able to change their behavior
 85 when the circumstances change. Serial reversal learning involves continuing to reverse the preference back
 86 and forth to determine whether individuals learn a “win-stay, lose-shift” rule that, when the reward is no
 87 longer in the expected option, they should switch to preferring the other option (Spence, 1936; J. Warren,
 88 1965; J. M. Warren, 1965). Once this rule is learned, it can then be applied to new contexts and result in
 89 improved performance over individuals who have not learned this rule (J. M. Warren, 1965). We randomly
 90 assigned individuals to a manipulated or control condition and used serial reversals (for the manipulated
 91 group) to attempt to manipulate flexibility and determine whether the manipulated individuals were then
 92 more flexible and more innovative in other contexts.

93 If grackle flexibility is manipulatable using serial reversals, this would provide us with a useful tool for in-
 94 vestigating the relationship between flexibility and any number of other variables implicated in geographic
 95 range expansions. It would provide researchers with a way to examine the direct links between, for example,
 96 flexibility and exploration, to determine whether they are connected and in which direction, which could
 97 provide insights into how populations establish in a new location if cross-population manipulations were con-
 98 ducted. If the flexibility manipulation is not successful, this could indicate either that we did not manipulate
 99 the right aspect of flexibility (e.g., perhaps training them to solve a variety of different types of tasks quickly
 100 would be more effective) or that grackle flexibility is not a trait that is trainable.



102 **Figure 1.** A visual illustration of Hypothesis 1 (A), Hypothesis 2 (B), and Hypothesis 4 (C). Longer black
 103 arrows indicate slower reversal times, the two yellow circles represent experience with the two yellow tubes
 104 that both contained food for the control group.



105

106 **Figure 2.** The experimental apparatuses: reversal learning using dark gray and light gray tubes or two
 107 different shapes on a touchscreen, and the wooden and plastic multi-access boxes (MAB). The wooden MAB
 108 has four loci, each containing food and each locus has a distinct way of being opened: lift up flap (A),
 109 swing open flap (B), pull out drawer (C), or push in flap (D). The plastic MAB has four loci that all provide access
 110 to one piece of food and each locus has a distinct way of being opened: open the window (left side), pull the
 111 string (top side), push the shovel (right side), or twist the shovel (bottom side).

112 **HYPOTHESES**

113 **H1: Behavioral flexibility, as measured by reversal learning using colored tubes, is manipulatable.**
 114

- 115 • **Prediction 1:** Individuals improve their flexibility on a serial reversal learning task using colored
 116 tubes by generally requiring fewer trials to reverse a preference as the number of reversals increases
 117 (manipulation condition). Their flexibility on this test is manipulated relative to control birds who do
 118 not undergo serial reversals. Instead, individuals in the control condition are matched to manipulated
 119 birds for experience (they experience a similar number of trials), but there is no possibility of a
 120 functional tube preference because both tubes are the same color (yellow) and both contain food,
 121 therefore either choice is correct.

122 • **P1 alternative 1:** If the number of trials to reverse a preference does not correlate with or positively
123 correlates with reversal number, which would account for all potential correlation outcomes, this sug-
124 gests that some individuals may prefer to rely on information acquired previously (i.e., they are slow
125 to reverse) rather than relying on current cues (e.g., the food is in a new location) (Griffin & Guez,
126 2014; Liu et al., 2016; e.g., Manrique et al., 2013; but see Homberg et al., 2007).

127 **H2: Manipulating behavioral flexibility (improving reversal learning speed through serial re-**
128 **versals using colored tubes) improves flexibility (rule learning and/or switching) and problem**
129 **solving in a new context (two distinct multi-access boxes and serial reversals on a touchscreen).**

130 • **P2:** Individuals that have improved their flexibility on a serial reversal learning task using colored
131 tubes (requiring fewer trials to reverse a preference as the number of reversals increases) are faster to
132 switch between new methods of solving (latency to solve or attempt to solve a new way of accessing
133 the food [locus]), and learn more new loci (higher total number of solved loci) on multi-access box
134 flexibility tasks, and are faster to reverse preferences in a serial reversal task using a touchscreen than
135 individuals in the control group where flexibility has not been manipulated. The positive correlation
136 between reversal learning performance using colored tubes and a touchscreen (faster birds have fewer
137 trials) and the multi-access boxes (faster birds have lower latencies) indicates that all three tests
138 measure the same ability even though the multi-access boxes require inventing new rules to solve new
139 loci (while potentially learning a rule about switching: “when an option becomes non-functional, try
140 a different option”) while reversal learning requires switching between two rules (“choose light gray”
141 or “choose dark gray”) or learning the rule to “switch when the previously rewarded option no longer
142 contains a reward”. Serial reversals eliminate the confounds of exploration, inhibition, and persistence
143 in explaining reversal learning speed because, after multiple reversals, what is being measured is the
144 ability to learn one or more rules. If the manipulation works, this indicates that flexibility can be
145 influenced by previous experience and might indicate that any individual has the potential to move
146 into new environments (see relevant hypotheses in preregistrations on [genetics](#) (R1) and [expansion](#)
147 (H1)).

148 • **P2 alternative 1:** If the manipulation does not work in that those individuals in the experimental
149 condition do not decrease their reversal speeds more than control individuals, then this experiment
150 elucidates whether general individual variation in flexibility relates to flexibility in new contexts (two
151 distinct multi-access boxes and serial reversals on a touchscreen) as well as problem solving ability
152 (multi-access boxes). The prediction is the same as in P2, but in this case variation in flexibility is
153 constrained by traits inherent to the individual (some of which will be tested in McCune KB et al.,
154 2019), which suggests that certain individuals will be more likely to move into new environments.

155 • **P2 alternative 2:** If there is no correlation between reversal learning speed (colored tubes) and the
156 latency to solve/attempt a new locus on the multi-access boxes, this could be because the latency
157 to solve not only measures flexibility but also innovativeness. In this case, an additional analysis is
158 run with the latency to solve as the response variable, to determine whether the fit of the model (as
159 determined by the lower AIC value) with reversal learning as an explanatory variable is improved if
160 motor diversity (the number of different motor actions used when attempting to solve the multi-access
161 box) is included as an explanatory variable (see Diquelou et al., 2015; Griffin et al., 2016). If the
162 inclusion of motor diversity improves the model fit, then this indicates that the latency to solve a
163 new locus on the multi-access box is influenced by flexibility (reversal learning speed) and innovation
164 (motor diversity).

165 • **P2 alternative 3:** If there is a negative correlation or no correlation between reversal learning speed
166 on colored tubes and reversal learning speed on the touchscreen, then this indicates that it may be
167 difficult for individuals to perceive and/or understand images on the touchscreen in contrast with
168 physical objects (colored tubes) (e.g., O’Hara et al., 2015).

169 **H4: Individuals should converge on an epsilon-first learning strategy (learn the correct choice**
170 **after one trial) as they progress through serial reversals.**

- 171 • **P4:** Individuals prefer a mixture of learning strategies in the first serial reversals (an *epsilon-decreasing*
172 strategy where individuals explore both options extensively before learning to prefer the rewarded op-
173 tion, and an *epsilon-first* strategy where the correct choice is consistently made after the first trial),
174 and then move toward the epsilon-first learning strategy. The epsilon-first strategy works better later
175 in the serial reversals where the reward is all or nothing because individuals have learned the environ-
176 ment is changing in predictable ways (Bergstrom & Lachmann, 2004): only one option is consistently
177 rewarded, and if the reward isn't in the previously rewarded option, it must be in the other option.
- 178 • **P4 alternative 1:** Individuals continue to prefer a mixture of learning strategies, and/or they do not
179 converge on the more functional epsilon-first learning strategy, regardless of how many reversals they
180 participate in. This pattern could suggest that the grackles do not attend to functional meta-strategies,
181 that is, they do not learn the overarching rule (once food is found in the non-preferred tube, one must
182 switch to preferring that tube color), but rather they learn each preference change as if it was new.

183 METHODS

184 Please see our preregistration that received in principle acceptance at PCI Ecology ([PDF](#) version) for all of
185 the preregistered methods. Below, we include a summary and describe all deviations from the preregistration.
186 We present the results from different hypotheses in separate articles: this one, K. McCune et al. (2022), and
187 Lukas et al. (2022).

188 **Planned Sample** Great-tailed grackles were caught in the wild in Tempe, Arizona, USA for individual
189 identification (colored leg bands in unique combinations). Some individuals (~32: ~16 in the control group
190 (they receive 1 reversal) and ~16 in the flexibility manipulation (they receive multiple reversals)) were brought
191 temporarily into aviaries for testing, and then released back to the wild.

- 192 • **Deviation from the plan:** we were able to test a total of 20 individuals: 11 in the control condition
193 and 9 in the manipulation condition. This met our minimum sample size criterion (see next section).

194 **Data collection stopping rule** We stopped testing birds after we completed two full aviary seasons
195 because the sample size was above the minimum suggested boundary of 15 (to detect a medium effect size)
196 based on model simulations (see Supplementary Material 6).

197 **Open materials** [Design files](#) for the plastic multi-access box: 3D printer files and laser cutter files

198 [Testing protocols](#) for all three experiments: colored tube reversal learning, plastic multi-access box, wooden
199 multi-access box, and touchscreen reversal learning

200 **Open data** The data are available at the Knowledge Network for Biocomplexity's data repository: https://knb.ecoinformatics.org/view/corina_logan.84.42.
201

202 **Randomization and counterbalancing** H1: Subjects were randomly assigned to the manipulated or
203 control group. In the reversal learning trials, the rewarded option is pseudorandomized for side (and the
204 option on the left is always placed first). Pseudorandomization consisted of alternating location for the first
205 two trials of a session and then keeping the same color on the same side for at most two consecutive trials
206 thereafter. A list of all 88 unique trial sequences for a 10-trial session, following the pseudorandomization
207 rules, was generated in advance for experimenters to use during testing (e.g., a randomized trial sequence

208 might look like: LRLRLRLRLR, where L and R refer to the location, left or right, of the rewarded tube).
209 Randomized trial sequences were assigned randomly to any given 10-trial session using a random number
210 generator (random.org) to generate a number from 1-88. The only exception to this randomization was when
211 an individual exhibited a side bias (choosing one side 4 or more trials in a row). In these cases, we stopped
212 the current random numbers for side and started putting the rewarded color on the non-preferred side as
213 much as possible while still following the pseudorandomization rules until the individual stopped exhibiting
214 a side bias.

215 ANALYSIS PLAN

216 Analyses were conducted in R [current version 4.1.2; R Core Team (2017)], using several R packages: Zhu
217 (2021), Hlavac (2018), Hadfield (2010), Bartoń (2020), McElreath (2020), Stan Development Team (2020),
218 Xie (2019), Ushey et al. (2020), Eddelbuettel & François (2011), Wickham (2016), knitr (Xie, 2013, 2017,
219 2018), Wickham et al. (2021), Gabry & Češnovar (2021), posterior (Bürkner et al., 2020), cowplot (Wilke,
220 n.d.), bayesplot (Gabry et al., 2019), irr (Gamer et al., 2012), psych (Revelle, 2014, 2017), Lin (2020),
221 DHARMA (Hartig, 2019), lme4 (Bates et al., 2012; Bates et al., 2015).

222 **Unregistered analyses:** We conducted unregistered interobserver reliability analyses on the response
223 variables. Scores indicated that the response variables are repeatable to a high or extremely high degree
224 given our instructions and training (see Supplementary Material 5).

225 **Planned analyses:** When there is more than one experimenter within a test, experimenter will be added as
226 a random effect to account for potential differences between experimenters in conducting the tests. If there
227 are no differences between models including or excluding experimenter as a random effect, then we will use
228 the model without this random effect for simplicity.

- 229 • **Deviation from the plan:** We removed experimenter (random variable) from all analyses because the
230 interobserver reliability scores were so high, indicating there was no difference between experimenters,
231 therefore we could keep our models simpler by leaving this variable out.

232 **Data checking** The data were checked for overdispersion, underdispersion, zero-inflation, and het-
233 eroscedasticity with the DHARMA R package (Hartig, 2019) following methods by [Hartig](#). Note: DHARMA
234 doesn't support MCMCglmm, therefore we will use the closest supported model: glmer from the R package
235 lme4 (Bates et al., 2015) for the DHARMA data checking.

236 **Determining the threshold: How many reversals are enough? The plan:** We initially (in 2017)
237 set as the passing criterion: During the data collection period, the number of trials required to reverse
238 a preference will be documented per bird, and reversals will continue until the first batch of birds tested
239 reaches an asymptote (i.e., there are negligible further decreases in the number of trials required to reverse a
240 preference). The number of reversals to reach the asymptote will be the number of reversals that subsequent
241 birds experience.

- 242 • **Deviations from the plan:** Due to delays in setting up the field site, we were only able to test two
243 grackles in early 2018 (January through April) and, due to randomization, only one (Fajita) was in the
244 experimental condition that involved undergoing the flexibility manipulation (Empanada was in the
245 control condition). While Fajita's reversal speeds generally improved with increasing serial reversals,
246 she never reached an asymptote (which we defined as passing three consecutive reversals in the same
247 number of trials), even after 38 reversals. These 38 reversals took 2.5 months, which is an impractical
248 amount of time if birds are to participate in the rest of the test battery (multi-access box, detour,
249 causal cognition, go no-go, reversal on a touchscreen) after undergoing the reversal manipulation (we
250 were initially permitted to keep them in aviaries for up to three months per bird, which we extended
251 to 6 months per bird in Dec 2018). Because our objective in this experiment was to manipulate an
252 individual's flexibility, we decided to revise our serial reversal passing criterion to something more

253 species relevant based on Fajita’s serial reversal performance and the performance of seven grackles in
254 Santa Barbara who underwent only one reversal in 2014 and 2015 (Logan, 2016). **The revised serial**
255 **reversal passing criterion was: passing two reversals in a row at or under 50 trials.** 50
256 trials is fewer trials than any of the nine grackles required to pass their first reversal (range 70-130),
257 therefore it should reflect an improvement in flexibility.

258 **Revising the choice criterion and the criterion to pass the control condition** **Choice criterion:**
259 At the beginning of the second bird’s initial discrimination in the reversal learning colored tube experiment
260 (October 2018), we revised the criterion for what counts as a choice from A) the bird’s head needs to pass
261 an invisible line on the table that ran perpendicular to the the tube opening to B) the bird needs to bend its
262 body or head down to look in the tube. Criterion A resulted in birds making more choices than the number
263 of learning opportunities they were exposed to (because they could not see whether there was food in the
264 tube unless they bent their head down to look in the tube) and appeared to result in slower learning. It is
265 important that one choice equals one learning opportunity, therefore we revised the choice criterion to the
266 latter. Anecdotally, this choice matters because the first three birds in the experiment (Tomatillo, Chalupa,
267 and Queso) learned faster than the pilot birds (Empanada and Fajita) in their initial discriminations and
268 first reversals. Thus, it was an important change to make at the beginning of the experiment (after testing
269 the two pilot birds and before collecting any data that were included in analyses).

270 **Criterion to pass the control condition:** Before collecting experimental data, we set the number of
271 trials experienced by the birds in the control group as 1100 because this is how many trials it would have
272 taken the pilot bird in the manipulated group, Fajita, to pass serial reversals 2-17 according to our revised
273 serial reversal passing criterion. However, after 25 and 17 days (after Tomatillo and Queso’s first reversals,
274 respectively) of testing the first two individuals in the control group, it became apparent that 1100 trials
275 is impractical given the time constraints for how long we were permitted to keep each bird temporarily in
276 captivity and would prevent birds from completing the test battery before their release. Additionally, after
277 revising the choice criterion, it was going to be likely that birds in the manipulated group would require fewer
278 than 1100 trials to meet the serial reversal passing criterion. Therefore, reducing the number of trials the
279 control birds experience would result in a better match of experience with birds in the manipulated group.
280 On 2 November 2018 we set the number of trials control birds experience after their first (and only) reversal
281 to the number of trials it requires the first bird in the manipulated group to pass (the first bird has not
282 passed yet, therefore we do not yet know what this number is). After more individuals in the manipulated
283 group passed, we updated this number to the average number of trials to pass. This applied to all birds in
284 the control condition, except Mofongo (see next paragraph).

- 285 • **Deviation from the plan** (16 April 2020): Mofongo (control condition) was a slow participator and
286 would not have finished his test battery by the time it got too hot to keep birds in the aviaries if we
287 used the current average number of trials (420). Instead, we matched him with the fastest bird in the
288 manipulated group (Habanero=290 trials) to make it more likely that Mofongo could get through the
289 rest of the test battery in time.

290 **P1: negative relationship between the number of trials to reverse a preference and the number**
291 **of reversals?** **Analysis:** Response variable: Number of trials to reverse a preference. An individual is
292 considered to have a preference if it chose the rewarded option at least 17 out of the most recent 20 trials
293 (with a minimum of 8 or 9 correct choices out of 10 on the two most recent sets of 10 trials). We use
294 a sliding window to look at the most recent 10 trials for a bird, regardless of when the testing sessions
295 occurred. Explanatory variable: reversal number. Random variables: batch (random effect because multiple
296 batches included in the analysis; batch is a test cohort, consisting of 8 birds being tested simultaneously)
297 and ID (random effect because repeated measures on the same individuals). A Generalized Linear Mixed
298 Model [GLMM; MCMCglmm function, MCMCglmm package; Hadfield (2010)] will be used with a Poisson
299 distribution and log link using 13,000 iterations with a thinning interval of 10, a burnin of 3,000, and
300 minimal priors ($V=1$, $\nu=0$) (Hadfield, 2014). We will ensure the GLMM shows acceptable convergence [lag

301 time autocorrelation values <0.01 ; Hadfield (2010)], and adjust parameters if necessary. We will determine
302 whether an independent variable had an effect or not using the Estimate in the full model.

303 We did not need a power analysis to estimate our ability to detect actual effects because, by definition,
304 the individuals that complete this experiment must get faster at reversing in order to be able to pass the
305 stopping criterion (two consecutive reversals in 50 trials or less). According to previous grackle data (from
306 the pilot birds, and from Santa Barbara Logan, 2016), the fastest grackle passed their first reversal in 70
307 trials, which means that passing our serial reversal stopping criterion would require them to have improved
308 their passing speed.

309 **P2: serial reversal improves rule switching and problem solving** **Unregistered analysis:** Be-
310 cause the wooden multi-access box was added after in principle recommendation, we conducted an unreg-
311 istered analysis to determine whether the plastic and wooden multi-access box results correlated with each
312 other, which would indicate that these tests are interchangeable. We found that they did not statistically
313 significantly correlate with each other on either variable measured: the average latency to attempt a new
314 locus (switching; Pearson's $r=0.74$, 89% confidence level= $0.02-0.95$, $t=2.18$, $df=4$, $p=0.09$, $n=6$) or the total
315 number of loci solved (problem solving; Pearson's $r=0.51$, 89% confidence level= $0.03-0.80$, $t=1.86$, $df=10$,
316 $p=0.09$, $n=12$). Therefore, while the performance on the two multi-access boxes might not be completely
317 independent as indicated by the high r values, the two boxes appear not to be completely interchangeable
318 either as indicated by the lack of statistical significance and high uncertainty in the r values. We therefore
319 analyzed the plastic and wooden multi-access boxes separately.

320 **Planned analyses:** As originally planned, we replaced the GLMs and GLMMs in May 2020 with more
321 powerful models after learning how to make bespoke Bayesian models from McElreath (2016). We made
322 these models before analyzing the actual data (14 May 2020).

323 One model was run per response variable: average latency to attempt to solve a new locus after solving
324 a different locus, and total number of loci solved. Explanatory variable: Number of trials to reverse a
325 preference in the last reversal. Random variable: batch.

326 The model for the number of loci solved takes the form of:

$$327 \text{locisolved} \sim \text{Binomial}(4, p) \text{ [likelihood]}$$

$$328 \text{logit}(p) \sim \alpha[\text{batch}] + \beta\text{trials} \text{ [model]}$$

329 locisolved is the number of loci solved on the multi-access box, 4 is the total number of loci on the multi-
330 access box, p is the probability of solving any one locus across the whole experiment, α is the intercept and
331 each batch gets its own, β is the expected amount of change in locisolved for every one unit change in trials,
332 and trials is the number of trials to reverse a color preference. See Supplementary Material 3 for more model
333 details.

334 The model for the latency to switch options takes the form of:

$$335 \text{latency} \sim \text{gamma-Poisson}(\lambda_i, \phi) \text{ [likelihood]}$$

$$336 \log(\lambda_i) \sim \alpha[\text{batch}] + \beta\text{trials} \text{ [model]}$$

337 latency is the average latency to attempt a new locus on the multi-access box, λ_i is the rate (probability of
338 attempting a locus in each second) per bird (and we take the log of it to make sure it is always positive; birds
339 with a higher rate have a smaller latency), ϕ is the dispersion of the rates across birds, α is the intercept
340 for the rate per batch, β is the expected amount of change in the rate of attempting to solve in any given
341 second for every one unit change in trials, and trials is the number of trials to reverse a color preference. See
342 Supplementary Material 6 for more model details.

343 **Deviations from the plan:**

- 344 • April 2020: we realized that the average latency to solve a new locus after solving a different locus is
345 confounded with the total number of loci solved because the measure of innovation is included in the

346 definition. Therefore, we removed average latency to solve a locus from analyses so that we are only
347 examining pure measures of flexibility (average latency to **attempt** to solve) and innovation (total
348 number of loci solved).

- 349 • Removed batch (random variable): the original model for P2 (Table SM2: Model 1) included the co-
350 variate aviary batch, however this ended up confounding the analysis because control and manipulated
351 individuals, while randomly assigned to these conditions, ended up in particular batches as a result
352 of their willingness to participate in tests offered during their time in the aviary (Table SM2: Model
353 3). Several grackles never passed habituation or training such that their first experiment could begin,
354 therefore we replaced these grackles in the aviaries with others who were willing to participate. This
355 means that batch did not indicate a particular temporal period. Therefore, we removed batch from
356 the models (post data collection, mid-data analysis).
- 357 • When making the bespoke Bayesian models, we realized that we had previously misinterpreted which
358 variable should be the response variable in this analysis. We originally set the number of trials to
359 reverse as the response variable, however we should have instead set the number of loci solved as the
360 response variable and then planned to conduct a second model with the latency to attempt a new
361 locus as the response variable and number of trials as the explanatory variable. This is because a) we
362 manipulated the number of trials to reverse, therefore it must be the explanatory variable (Hernán &
363 Robins, 2006); and b) they should be split into two models because of a and because these are two very
364 different relationships that should be considered in their own models. We also realized that Condition
365 (manipulated or control) does not need to be a variable in any of our models because the manipulated
366 birds have, by definition, faster reversal speeds.

367 **P2 alternative 2: additional analysis: latency and motor diversity** **Planned analysis:** We ran one
368 model per response variable: Number of trials to attempt a new locus on the multi-access boxes, and number
369 of trials to solve (meet criterion) a new locus on the multi-access boxes. Explanatory variables: Number of
370 trials to reverse a preference in the last reversal that individual participated in, motor diversity: the number
371 of different motor actions used when attempting to solve the multi-access boxes. Random variable: ID
372 (random because repeated measures on the same individuals). A Generalized Linear Mixed Model [GLMM;
373 MCMCglmm function, MCMCglmm package; Hadfield (2010)] will be used with a Poisson distribution and
374 log link using 13,000 iterations with a thinning interval of 10, a burnin of 3,000, and minimal priors ($V=1$,
375 $\nu=0$) (Hadfield, 2014). We ensured the GLMM showed acceptable convergence [lag time autocorrelation
376 values <0.01 ; Hadfield (2010)] by adjusting parameters if necessary. We determined whether an independent
377 variable had an effect or not using the Estimate in the full model.

- 378 • **Deviations from the plan:** We used the average latency rather than the number of trials to attempt
379 a new locus because this would make the model comparable with the model in P2. Using the number
380 of trials was an artifact from a previous version and we had missed updating this. We omitted the
381 number of trials to solve a new locus as described in the deviation from the plan in P2 above. We used
382 a GLM rather than a GLMM because there was only one data point per bird.

383 **P4: learning strategies (for birds in the manipulated group only)** **Analysis 1 (qualitative):**
384 Learning strategies were identified by matching them to the two known approximate strategies of the contex-
385 tual, binary multi-armed bandit: epsilon-first and epsilon-decreasing (McInerney, 2010; as in Logan, 2016).
386 We used the criterion for the epsilon-first strategy of learning the correct choice after one trial and then
387 choosing correctly thereafter. Other patterns were classified as the epsilon-decreasing strategy. This method
388 of qualitative inspection of learning curves is standard for this type of learning strategy assessment (McIn-
389 erney, 2010). The variable for visual inspection was the proportion of correct choices in a non-overlapping
390 sliding window of 4-trial bins across the total number of trials required to reach the criterion of 17/20 correct
391 choices per individual.

392 From Logan (2016) (emphasis added):

393 The following equations refer to the different phases involved in each strategy:

394 Equation 1 (exploration phase):

$$\epsilon N$$

395 Equation 2 (exploitation phase):

$$(1 - \epsilon)N$$

396 N is the number of trials given, and epsilon,

$$\epsilon$$

397 , represents the subject's uncertainty about the location of the reward, starting at complete
398 uncertainty ($\epsilon = 1$) at the beginning of the experiment and decreasing rapidly as individuals gain
399 experience with the task (exploration phase where the rewarded [option] is chosen below or at
400 chance levels) and switch to the exploitative phase (the rewarded [option] is chosen significantly
401 above chance levels). Because the [subjects] needed to learn the rules of the task, they necessarily
402 had an exploration phase. The **epsilon-first strategy** involves an exploration phase followed
403 by an entirely exploitative phase. The optimal strategy overall would be to explore one color in
404 the first trial and the other color in the second trial, and then switch to an exploitative strategy
405 (choose the rewarded [option] significantly above chance levels). In this case there would be
406 no pattern [in the learning curve] in the choices [during] the exploration phase because it would
407 consist of sampling each [option] only once. In the **epsilon-decreasing strategy**, subjects would
408 start by making some incorrect choices and then increase their choice of the rewarded [option]
409 gradually as their uncertainty decreases until they choose the rewarded [option] significantly
410 above chance levels. In this case, a linear pattern emerges (in the learning curve) during the
411 exploration phase.

412 **Analysis 2 (quantitative):** We then quantitatively determined to what degree each bird used the explo-
413 ration versus exploitation strategy using methods in (Federspiel et al., 2017) by calculating the number of
414 20-trial blocks where birds were choosing “randomly” (6-14 correct choices; called sampling blocks; akin to
415 the exploration phase above) and dividing it by the total number of blocks to reach criterion per bird. This
416 ratio was also calculated for “acquisition” blocks where birds made primarily correct choices (15-20 correct
417 choices; akin to the exploitation phase above). These ratios, calculated for each bird for their serial reversals,
418 quantitatively discern the exploration from the exploitation phases.

- 419 • **Deviation from the plan** (Aug 2021): the grackles were tested in 10-trial blocks and not 20-trial
420 blocks as in Federspiel et al. (2017), which would mean that if there were <20 trials in the last block
421 of a reversal, they would be omitted from the analysis. Therefore, we changed the block size to 10
422 trials and adjusted the sampling blocks to 2-9 correct choices, and the acquisition blocks to 9-10 correct
423 choices using significance levels in the binomial test as did Federspiel et al. (2017).

424 DEVIATIONS FROM THE PREREGISTRATION

425 In the middle of data collection

- 426 1) 10 April 2019: We **discontinued the reversal learning experiment on the touchscreen** because
427 it appears to measure something other than what we intended to test and it requires a huge time
428 investment for each bird (which consequently reduces the number of other tests they are available
429 to participate in). This is not necessarily surprising because this is the first time touchscreen tests
430 have been conducted in this species, and also the first time (to our knowledge) this particular reversal
431 experiment has been conducted on a touchscreen with birds. We based this decision on data from four
432 grackles (2 in the flexibility manipulation group and 2 in the flexibility control group; 3 males and 1
433 female). All four of these individuals showed highly inconsistent learning curves and required hundreds
434 more trials to form each preference when compared to the performance of these individuals on the

435 colored tube reversal experiment. It appears that there is a confounding variable with the touchscreen
436 such that they are extremely slow to learn a preference as indicated by passing our criterion of 17 correct
437 trials out of the most recent 20. We will not include the data from this experiment when conducting
438 the cross-test comparisons in the Analysis Plan section of the preregistration. Instead, in the Results
439 section, we provide summary results for this experiment and, in the Discussion, qualitatively compare
440 it with performance on the colored tube reversal test to explain what might have confounded the
441 touchscreen experiment.

- 442 2) 16 April 2019: Because we discontinued the touchscreen reversal learning experiment, we **added an**
443 **additional but distinct multi-access box** task, which allowed us to continue to measure flexibility
444 across three different experiments. There are two main differences between the first multi-access box,
445 which is made of plastic, and the new multi-access box, which is made of wood. First, the wooden
446 multi-access box is a natural log in which we carved out 4 compartments. As a result, the apparatus and
447 solving options are more comparable to what grackles experience in the wild, though each compartment
448 is covered by a transparent plastic door that requires different behaviors to open. Furthermore, there
449 is only one food item available in the plastic multi-access box and the bird could use any of 4 loci
450 to reach it. In contrast, the wooden multi-access box has a piece of food in each of the 4 separate
451 compartments.

452 Post data collection, pre-data analysis

- 453 3) We completed our simulation to explore the lower boundary of a minimum sample size and determined
454 that **our sample size for the Arizona study site is above the minimum** (see details and code
455 in Supplementary Material 1; 17 April 2020).
- 456 4) Please see our Alternative Analyses section in the preregistration where we stated that we would
457 learn and implement Bayesian models, which resulted in our **changing the analysis for P2** and
458 that we are replacing this analysis with the new models in the Ability to detect actual effects section
459 (Supplementary Material 1; 14 May 2020). We also describe in SM1 that we realized that Condition
460 (manipulated or control) does not need to be a variable in our models because the manipulated birds
461 have, by definition, faster reversal speeds.
- 462 5) We originally planned on testing only **adults** to have a better understanding of what the species is
463 capable of, assuming the abilities we are testing are at their optimal levels in adulthood, and so we
464 could increase our statistical power by eliminating the need to include age as an independent variable
465 in the models. Because the grackles in Arizona were extremely difficult to catch, we ended up testing
466 two juveniles: Taco and Chilaquile. We did not conduct the full test battery with Taco or put him in
467 the flexibility manipulation or control groups (he received 1 reversal and then moved on to the next
468 test) because he was the first juvenile and we wanted to see whether his performance was different
469 from adult performances. His performances were similar to the adults, therefore we decided to put
470 Chilaquile in the full test battery. Chilaquile's performances were also similar to the adults, therefore
471 we decided not to add age as an independent variable in the models to avoid reducing our statistical
472 power.

473 Post data collection, mid-data analysis

- 474 6) The original model for P2 (Table SM3: Model 1) included the covariate aviary batch, however this ended
475 up confounding the analysis because control and manipulated individuals, while randomly assigned to
476 these conditions, ended up in particular batches as a result of their willingness to participate in tests
477 offered during their time in the aviary (Table SM3: Model 3). Several grackles never passed habituation
478 or training such that their first experiment could begin, therefore we replaced these grackles in the
479 aviaries with others who were willing to participate. This means that batch did not indicate a particular
480 temporal period. Therefore, we **removed batch from the model**.

481 **RESULTS**

482 Data are publicly [available](#) at the Knowledge Network for Biocomplexity (C. Logan et al., 2022). Although
 483 22 grackles completed their initial colored tube discrimination, only 20 grackles participated in one or more
 484 reversals (Table SM5). The rest of the tests began only after a bird’s reversal experiment was complete (C.
 485 Logan et al., 2022).

486 **P1: reversal speed gets faster with serial reversals**

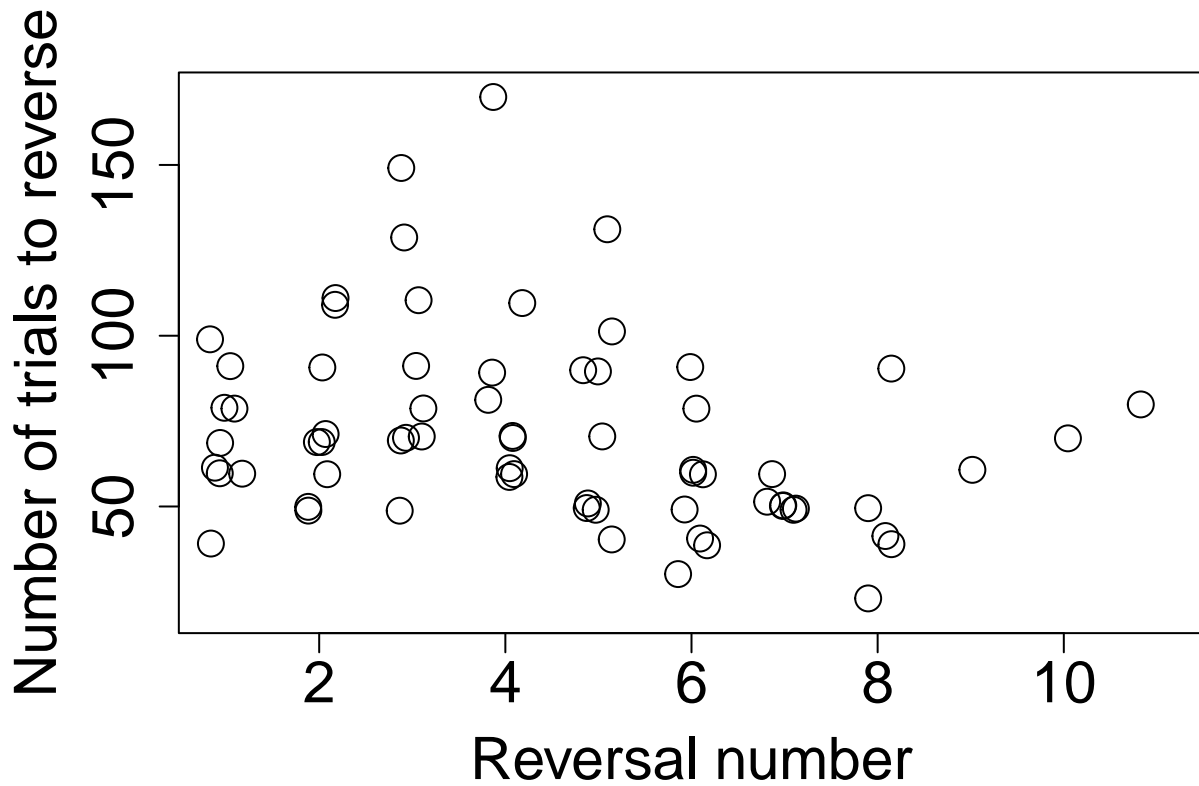
487 The birds in the manipulated group required a similar number of trials during their first reversal (R1 me-
 488 dian=75 trials) as the birds in the control group needed during their first and only reversal (R1 median=70
 489 trials) (see unregistered analysis in Table 1). The manipulated birds improved during the reversal manip-
 490 ulation to a median of 40 trials in their last reversal: there was a significant negative correlation between
 491 the number of trials to reverse (average=71 trials, standard deviation (sd)=28, Table 2) and the reversal
 492 number for those grackles in the flexibility manipulation condition (n=9, which included Memela who did
 493 not pass the manipulation condition of passing two consecutive reversals in 50 trials or less; Figure 3).

494 **Table 1.** Unregistered analysis: the number of trials to reverse in the first reversal is similar between the
 495 manipulated and control groups.

	Posterior mean	Lower 89 percentile compatibility interval (5.5%)	Upper 89 percentile compatibility interval (94.5%)	Effective sample size	pMCMC	Significance code: **=0.01
Intercept	4.29	4.12	4.46	420	<0.002	**
Manipulation Condition	-0.08	-0.27	0.11	420	0.46	

498 **Table 2.** The number of trials to reverse decreases with increasing reversal number.

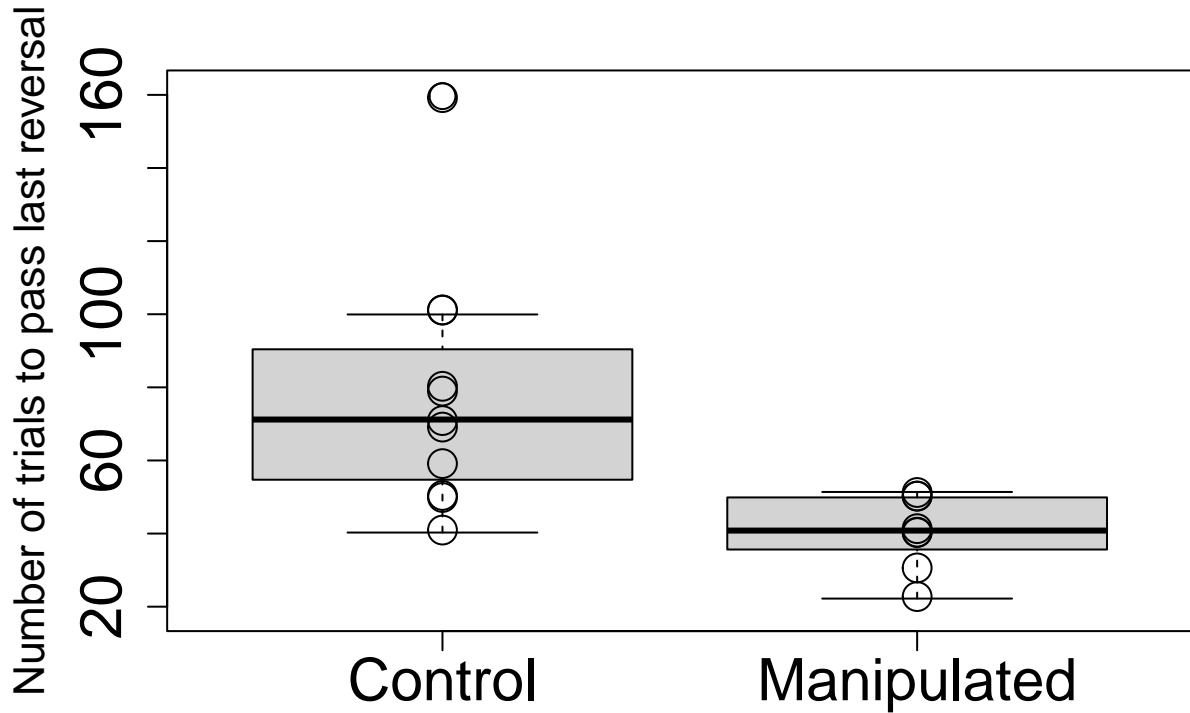
	Posterior mean	Lower 89 percentile compatibility interval (5.5%)	Upper 89 percentile compatibility interval (94.5%)	Effective sample size	pMCMC	Significance code: **=0.01
Intercept	4.44	4.31	4.62	420	<0.002	**
Reverse Number	-0.06	-0.10	-0.03	420	<0.002	**



501

502 **Figure 3.** Individuals in the manipulated condition (who received serial reversals) did not linearly decrease
 503 their reversal passing speeds with increasing reversal number (n=9 grackles).

504 **Unregistered analysis 1:** There was additionally a difference between manipulated and control reversal
 505 speeds when comparing their last reversals (Figure 4; for the control birds, their last reversal was their first
 506 reversal; Table 3). This analysis includes 19 grackles (8 manipulated condition - only those who actually
 507 passed the manipulation, 11 control condition) who had an overall average of 62 trials in their last reversal
 508 (sd=32).



509

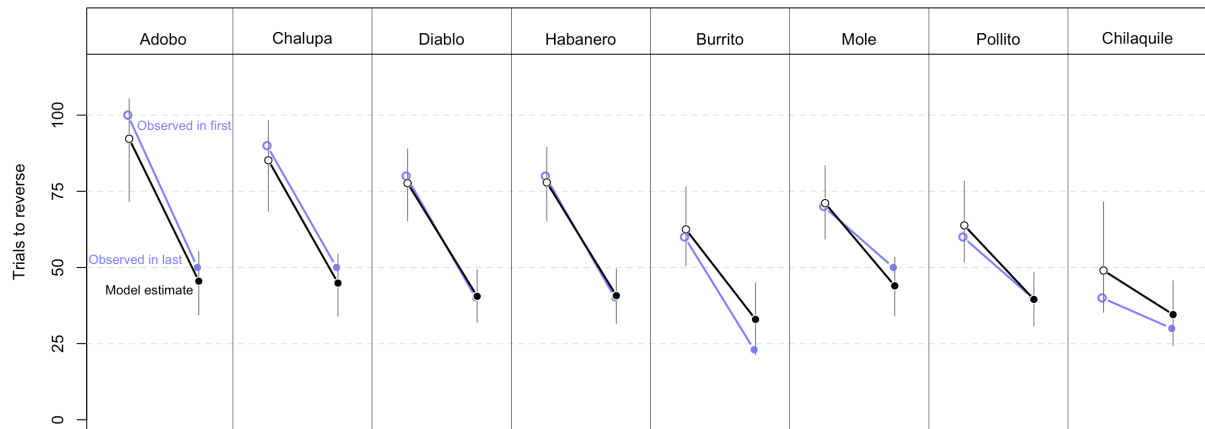
510 **Figure 4.** Individuals in the manipulated condition (who received serial reversals) passed their last reversal
 511 in fewer trials than individuals in the control condition (who only received 1 reversal). n=19 grackles:
 512 11=control, 8=manipulated.

513 **Table 3.** Individuals in the manipulated condition pass their last reversal in fewer trials than control
 514 individuals.

	Posterior mean	Lower 89 percentile compatability interval (5.5%)	Upper 89 percentile compatability interval (94.5%)	Effective sample size	pMCMC	Significance code: **=0.01
Intercept	4.28	4.08	4.48	420	<0.002	**
Reverse Number	-0.51	-0.81	-0.22	420	0.010	**

516

517 **Unregistered analysis 2:** A pooled model of performance across all reversals estimates that birds can
 518 expect to improve by about 30 trials (89% percentile interval (PI): 25-36; Table 7: Model 15) after completing
 519 the serial reversals. While all manipulated birds improved, those birds that were already fast to reverse in
 520 their first reversal improved less than the birds that required many trials to reverse in their first reversal
 521 (posterior peak indicates a correlation of +0.64, with highest posterior density intervals (HPDI) all positive,
 522 between the first reversal value and the improvement achieved by the last reversal; Table SM3: Model 16).
 523 However, the birds who were the fastest in the first reversal, were also the fastest in the last reversal, but
 524 the difference between the slower and faster reversers is reduced (Figure 5).


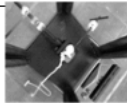





525

526 **Figure 5.** All eight manipulated birds needed fewer trials to reverse in their last reversal than in their
 527 first. Their improvement depended on their starting value, with steeper slopes for those birds that needed
 528 more trials to reverse in the first reversal (blue = observed values and changes, black = model estimates).
 529 However, birds who needed more trials in the first reversal did not completely catch up, such that the birds
 530 that needed more trials in their first reversal also needed more trials in their last reversal relative to other
 531 grackles.

532 **P2: serial reversals improve rule switching and problem solving on the MAB**

533 To determine whether the serial reversal manipulation affected flexibility generally, we compared performance
 534 (the number of trials to reverse a preference in the first and last color reversal, performance of the manipulated
 535 group relative to the control group) to speed of solution switching on two multi-access boxes. Furthermore,
 536 we assessed whether flexibility measured through these serial reversals related to innovativeness by comparing
 537 performance to the number of loci solved on the multi-access boxes. The results for each of these comparisons
 538 are described in detail below and an overview is provided in Figure 6.

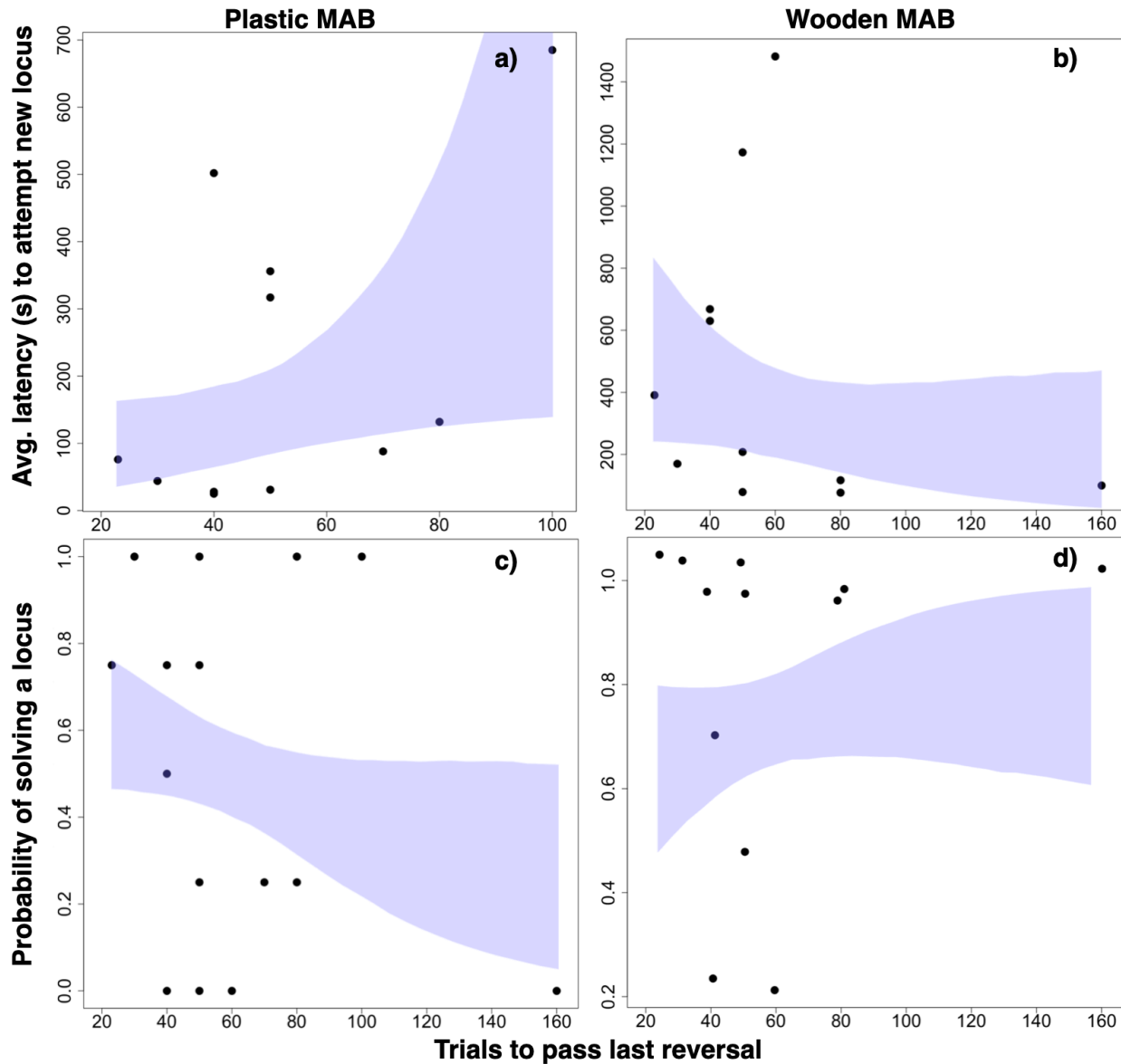
P2: How does flexibility, measured via performance on serial reversals, relate to flexibility in another context and innovativeness?		Flexibility (number of trials to pass in serial reversals) 		
		First Reversal	Last Reversal	Manipulated relative to Control
Flexibility in a new context (latency to switch loci)		+	+	+
		-	0	0
Innovativeness (number of loci solved)		0	+	0
		0	0*	+

539

540 **Figure 6.** Overview of the results from the P2 analyses with the multi-access boxes (plastic and wooden).
 541 An effect of natural variation in flexibility on performance on the multi-access box tasks would result in
 542 correlations in the first reversal. An effect of the flexibility manipulation would result in a change in cor-
 543 relations from the first to last reversals. A plus sign (+) indicates a positive correlation, a minus sign (-)
 544 indicates a negative correlation, and a 0 indicates no correlation between the two variables. The asterisks
 545 (*) indicate that a small sample size decreases the reliability of this result.

546 **Rule switching: latency to attempt a new locus on the multi-access box (plastic) ~ trials to**
 547 **reverse**

548 Grackles that were faster to reverse a preference in their **last reversal** (average 52 trials, sd=23), where
 549 grackles in the control condition received only one reversal which served as their first and last reversal, were
 550 also faster to attempt to solve a new locus on the plastic multi-access box (after just having passed criterion
 551 on a different locus; average=208 seconds, sd=226; Figure 7a; Table SM3: Model 9; n=11 grackles: 6 in
 552 manipulated condition, 5 in control condition; 6 subjects completed this experiment but solved 0 loci or 1
 553 locus and so did not have switching times). We also found that individuals in the flexibility manipulation
 554 had faster switch latencies than those in the control condition (Table SM3: Model 10). There was a positive
 555 correlation between the number of trials to reverse in the **first reversal** (average=70 trials, sd=21) and the
 556 average switch latency on the plastic multi-access box (Table SM3: Model 11). A correlation was determined
 557 to be present if the compatibility interval for the slope (b) in the model output did not cross zero (Table
 558 SM3). This criterion was used throughout the analyses for P2.



559

560 **Figure 7.** The average latency (seconds) to attempt to solve a different locus after having previously
 561 successfully solved a locus on a) the plastic multi-access box (MAB) is positively correlated with the number
 562 of trials to pass their last reversal (n = 11 grackles), but on b) the wooden MAB it is not correlated with
 563 the number of trials to pass their last reversal (n = 11 grackles). Additionally, the probability of solving a
 564 locus on c) the plastic MAB is negatively correlated with the number of trials to pass their last reversal (n
 565 = 15 grackles), but on d) the wooden MAB it is not correlated with the number of trials to pass their last
 566 reversal (n = 12 grackles, estimate of slope includes zero). Shading represents the 89 percentile compatibility
 567 intervals.

568 **Rule switching: latency to attempt a new locus on the multi-access box (wooden) ~ trials to**
 569 **reverse (unregistered analysis)**

570 There was no correlation between the number of trials to reverse a preference in their **last reversal** (average
 571 60 trials, sd=38) and the latency to attempt to solve a new locus on the wooden multi-access box (after just
 572 having passed criterion on a different locus; average=463 seconds, sd=481; Figure 7b; Table SM3: Model 12;
 573 n=11 grackles: 5 in manipulated condition, 6 in control condition; Diablo also completed this experiment
 574 and solved 1 locus, but did not attempt another locus after that, thus he does not have any switching times

575 to analyze). We additionally found that there was no difference in the average latency to switch between
 576 individuals in the flexibility manipulation and those in the control condition (Table SM3: Model 13). There
 577 was a negative correlation between the number of trials to reverse in the **first reversal** (average=73 trials,
 578 sd=34) and the average switch latency on the multi-access box (Table SM3: Model 14).

579 **Innovativeness: number of loci solved on the multi-access box (plastic) ~ trials to reverse**
 580 Grackles that were faster to reverse a preference in their **last reversal** (average 62 trials, sd=34) solved
 581 more loci on the plastic multi-access box (average=2 loci, sd=1.6; Figure 7c; Table SM3: Model 2; n=15
 582 grackles: 6 in manipulated condition, 9 in control condition; this number excludes Mole and Habanero who
 583 were, due to experimenter error, given the fully put together box during habituation and could have learned
 584 how to solve the loci at that time). There was no correlation between the number of loci solved and which
 585 reversal condition a grackle was randomly assigned to (Table SM3: Model 4). There was also no correlation
 586 between the number of trials to reverse in the **first reversal** (average=75 trials, sd=31) and the number of
 587 loci solved on the multi-access box (Table SM3: Model 5).

588 **Innovativeness: number of loci solved on the multi-access box (wooden) ~ trials to reverse**
 589 **(unregistered analysis)**

590 The compatibility interval for the estimate for the association (mean beta -0.41) between the number of
 591 loci solved on the wooden multi-access box (average=3.2, sd=1.3) and the number of trials to reverse a
 592 preference in their **last reversal** (average=59 trials, sd=38) crossed zero (Figure 5d; Model 6, Table SM3;
 593 n=12 grackles: 6 in manipulated condition, 6 in control condition). This could mean that there is no
 594 association, however simulations in Supplementary Material 1 showed that we would not be able to reliably
 595 distinguish whether a small effect is different from zero with our sample size (with a simulated beta of -1 and
 596 an sd in the number of trials >10, the compatibility interval of the estimate crossed zero in all simulations ;
 597 Table SM1.2). We did find a correlation between the number of loci solved and which reversal condition a
 598 grackle was randomly assigned to, indicating the reversal manipulation appears to have affected performance
 599 on the wooden multi-access box. The model estimates that manipulated birds solved on average 1.2 more loci
 600 than birds in the control condition (Table SM3: Model 7, wooden; 89% compatibility intervals=0.34-2.14;
 601 n=12 grackles: 6 in manipulated condition, 6 in control condition). However, there is no association between
 602 the number of trials to reverse in the first reversal (average=74 trials, sd=34) and the number of loci solved
 603 on the multi-access box (Table SM3: Model 8, wooden).

604 **P2 alternative 2 (additional analysis): latency and motor diversity**

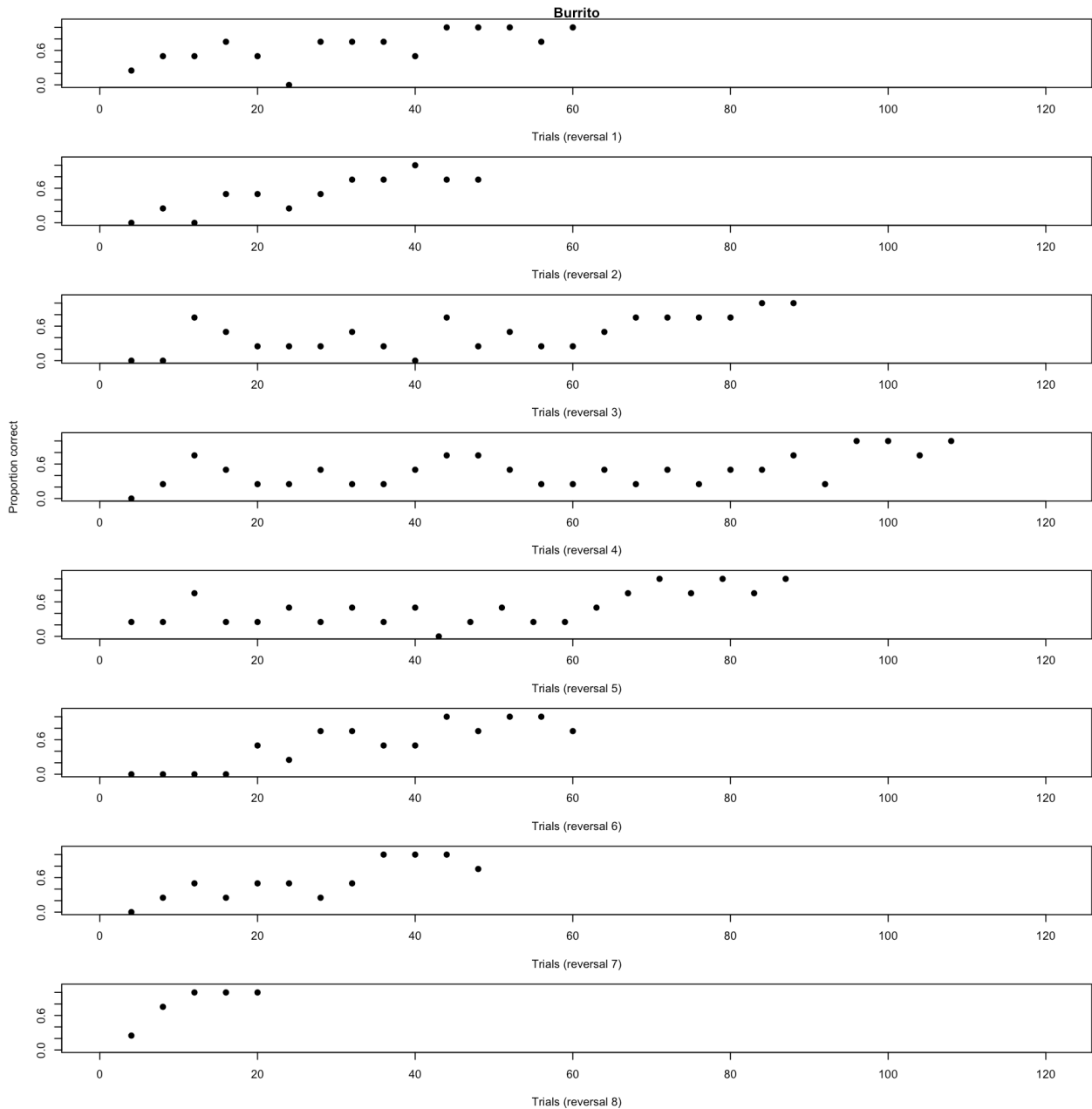
605 Because there was no correlation between the number of trials to reverse in the last reversal and the latency to
 606 attempt a different locus on the wooden multi-access box, we conducted this additional analysis to determine
 607 whether the model fit was improved when adding the number of motor actions as an explanatory variable.
 608 Adding the number of motor actions (wooden: average=13, sd=4) did not improve the model fit when
 609 examining the relationship between the latency to switch loci on the wooden multi-access box (wooden:
 610 average=463, sd=481) and the number of trials to reverse in the last reversal (wooden: average=60, sd=38)
 611 because the Akaike weights were similar for both models (wooden: n=11 grackles: 5 in the manipulated
 612 group, 6 in the control group; Table 4).

613 **Table 4.** Adding the number of motor actions used to the analysis of the average latency to attempt a
 614 new option on the wooden multi-access box and the number of trials to reverse in the last reversal does not
 615 improve the model fit.

	Intercept	Motor actions (wooden)	Trials last reversal	df	log likelihood	AICc	delta	weight
616	463.2	NA	NA	2	-83.025	171.6	0.00	0.674
	934.6	-35.28	NA	3	-82.477	174.4	2.83	0.164
	665.8	NA	-3.362	3	-82.631	174.7	3.14	0.140
617	1250.0	-40.68	-4.040	4	-81.850	178.4	6.82	0.022

618 **P4: serial reversal learning strategy**

619 **Analysis 1 (qualitative):** Using the criterion for the epsilon-first strategy of learning the correct choice
620 after one trial and then choosing correctly thereafter, no grackle in this study used this strategy in any
621 reversal. All grackles used an epsilon-decreasing strategy in all reversals (Figure 8 and Supplementary
622 Material 6). We use Burrito's figures to illustrate the epsilon-decreasing strategy (Figure 8): the proportion
623 of trials he gets correct wanders up and down (epsilon-decreasing) until an asymptote at 0.8 is reached and
624 held.



625

626 **Figure 8.** Burrito's proportion of trials correct by trial number and reversal showing the epsilon-decreasing
627 learning strategy where options are explored before forming a preference.

628 **Analysis 2 (quantitative):** We additionally quantitatively determined to what degree each bird used
629 the exploration versus exploitation strategy using methods in Federspiel et al. (2017) by calculating the
630 number of 10-trial blocks where birds were choosing "randomly" (2-9 correct choices; called sampling blocks;

631 akin to the exploration strategy) divided by the total number of blocks to reach criterion per bird. This
632 ratio was also calculated for “acquisition” blocks where birds made primarily correct choices (9-10 correct
633 choices; akin to the exploitation strategy). There was no correlation between exploration (sampling ratio) or
634 exploitation (acquisition ratio) and reversal number (sampling: reversal estimate=-0.09, SE=0.11, z=-0.86,
635 p=0.39; acquisition: reversal estimate=0.00, SE=0.00, z=-0, p=1.00), indicating that the grackles did not
636 use a particular strategy earlier or later in their serial reversals.

637 DISCUSSION

638 We conducted a controlled experiment to evaluate whether serial reversal learning affected flexibility and
639 innovativeness in new contexts. We found that the number of trials to reverse decreased with increasing
640 reversal number, and, when examining last reversals, there was a difference between the manipulated and
641 control groups. This indicates that the flexibility manipulation was effective in that it manipulated reversal
642 learning speeds, suggesting that these individuals shifted toward a “win-stay, lose-shift” rule to learn to
643 reverse faster after more experience with reversing (Spence, 1936; J. Warren, 1965; J. M. Warren, 1965).
644 The manipulated individuals who increased their reversal learning speed, were then apparently able to apply
645 this to a new context, which resulted in better performance when compared with control individuals who
646 did not have the opportunity to learn. Previous research has also exploited the fact that most individuals
647 can learn to learn and have used serial reversals to show that such experience usually improves performance
648 when transferring to reversals involving different stimuli (e.g., visual vs. spatial, visual vs. visual in a new
649 combination) (Rayburn-Reeves et al., 2013; Schusterman, 1962; J. Warren, 1965, 1966).

650 While performance differed between the two multi-access boxes, the serial reversal flexibility manipulation did
651 affect flexibility in a new context, as well as innovativeness. Grackles that were faster to reverse a preference
652 in their first and last reversals, and those in the manipulated condition, were also faster to attempt to solve
653 a new locus on the **plastic** multi-access box. Similarly, the flexibility manipulation affected innovativeness
654 because grackles in the manipulated condition solved on average 1.2 more loci on the **wooden** multi-access
655 box than those birds in the control condition and there was a positive correlation between the number of
656 loci solved on the **plastic** multi-access box and the number of trials to reverse in the last reversal. That our
657 results were not consistent across first reversal, last reversal, and condition (Figure 4) on the two different
658 multi-access boxes could be due to the small sample sizes because even in the control group there were several
659 individuals who solved their first and only reversal in very few trials. Furthermore, the lack of correlation
660 between the number of trials to reverse in the first reversal and the number of loci solved on either multi-
661 access box indicates that flexibility is not an inherently utilized tool, but one that is shaped by experience.
662 If it was an inherently utilized tool, the variation in the number of trials to complete first reversals would
663 likely have resulted in a correlation with the number of loci solved.

664 Our results are in contrast with previous research on the correlation between flexibility performance, using
665 serial reversals, and innovation: Indian mynas that were faster to reverse, were slower to innovate (Griffin et
666 al., 2013). However, the Griffin et al. (2013) investigation was designed to evaluate the correlation between
667 the variables and not whether manipulating flexibility using serial reversals influenced innovativeness. This
668 difference could explain the differing results because correlational research can become noisy if there are
669 unmeasured variables, which is something that a manipulation can help reduce. Other potential reasons for
670 the difference in results could be due to using different experimental designs, and/or different serial reversal
671 passing criteria (Griffin et al., 2013 used a preset number of reversals that resulted in a maximum of four
672 reversals).

673 None of the flexibility manipulated individuals converged on using an epsilon-first learning strategy (learn
674 the correct choice after one trial) as they progressed through serial reversals. All used the epsilon-decreasing
675 strategy (explore options before forming a preference) throughout their reversals. Additionally, no grackle
676 used a particular exploitation or exploration strategy earlier or later in their reversals. Learning theory on
677 serial reversal experiments predicts that all individuals in the manipulated group shifted toward the “win-
678 stay, lose-shift” rule because their reversal speeds improved (Spence, 1936; J. Warren, 1965; J. M. Warren,
679 1965). In contrast, learning theory on multi-armed bandit (a paradigm often used in reversal learning)
680 decision making has a stricter criterion, predicting that the optimal strategy is to maximize the cumulative

681 reward, which, in this case would result in individuals using the epsilon-first learning strategy immediately
682 after the first trial (McInerney, 2010). Both learning theories consider one trial learning the optimal solution.
683 Perhaps these wild-caught grackles relied solely on the epsilon-decreasing strategy because these individuals
684 are used to an environment where information about the probability of what the optimal options are varies
685 (McInerney, 2010). Therefore, maximizing information gain via continued exploration of the available options
686 is likely of more use in the less predictable environment in the wild. Other investigations of the exploitation
687 vs. exploration learning strategies involved in reversal learning have found that these strategies can vary by
688 individual and relate to differences in reversal performance. For example, urban common mynas were slower
689 to reverse a preference than rural mynas because they spent more time exploring their options (Federspiel et
690 al., 2017). Perhaps we found no such differences in the grackles because all of the individuals we tested came
691 from an urban area. If a rural population of grackles could be found, it would be interesting to compare
692 learning strategy use between rural and urban individuals.

693 **Why did performance on a touchscreen vary so drastically from a traditional approach?**

694 We assumed that reversal learning performance using **shape on the touchscreen** would directly compare
695 to and be interchangeable with reversal learning performance using colored tubes. However, it quickly
696 became clear that the touchscreen experiment may have been asking a different question compared with
697 the traditional reversal learning approach using physical objects. Unfortunately, we did not have the time
698 to explore what might have caused the differences between the two tests, but we speculate below. We
699 conclude that these two methods, the traditional physical object and the touchscreen, do not measure the
700 same construct in this species and with this reversal learning experiment.

701 One possible explanation for the difference between the two experiments is that grackles might require more
702 trials to learn to discriminate between shapes than between colors. Shapes are known to require a few more
703 trials for a preference to develop (e.g., Shaw et al., 2015: mean=40 trials color, mean=55 trials shape in
704 toutouwai; Isden et al., 2013: mean=6 trials color, mean=10 trials shape in spotted bowerbirds), however
705 grackles required hundreds more trials to learn shapes, therefore this explanation seems unlikely. Moreover,
706 grackles may not have understood how the touchscreen worked and therefore it was the apparatus that
707 interfered with their performance, yet grackles successfully completed a go no-go inhibition task using the
708 same touchscreen apparatus (Logan et al., 2021). The go no-go task similarly used two different white
709 shapes (wavy lines or a heart), but the shapes were presented sequentially rather than simultaneously (as
710 in the reversal touchscreen experiment). Given this difference between the two touchscreen experiments, it
711 is possible that the grackles found touching the screen in the reversal experiment rewarding in and of itself
712 because something happened whenever they made a response. That is, if they touched the correct stimulus,
713 they received food; if they touched the incorrect stimulus, the screen went blank immediately. This is in
714 contrast with the go no-go experiment where the stimulus stayed on the screen for a set amount of time after
715 an incorrect choice. Another potential reason for the difference between performances on the two touchscreen
716 experiments was that making the incorrect choice in the reversal experiment was not costly enough. In the
717 reversal touchscreen experiment, they could get through many trials, receiving some rewards, in a short
718 amount of time. Consequently, there was potentially not enough incentive to learn quickly, thus explaining
719 the differences in learning speeds between the two reversal experiments.

720 We are not the first group to attempt to transfer a traditional lab or field task to a touchscreen apparatus
721 (e.g., Drayton & Santos, 2014). Despite some of the challenges associated with touchscreen apparatuses,
722 other attempts to transfer tasks to a touchscreen have been more successful (e.g., Blaisdell & Cook, 2005;
723 Kangas & Bergman, 2017; Sawa et al., 2005). We maintain that touchscreens have the potential to be an
724 incredibly useful tool for studying comparative cognition in some systems (for reviews and methods, see
725 Bussey et al., 2008; Cook et al., 2004; Kangas & Bergman, 2017; Logan et al., 2021; Seitz et al., 2021; Wolf
726 et al., 2014).

727 Conclusion

728 We demonstrate that it is possible to manipulate flexibility, using a paradigm such as reversal learning, to
729 examine its direct link with other traits. This opens up many opportunities for future research to better
730 understand what flexibility is and whether and how it is causally related to other behaviors or forms of
731 cognition. Understanding how flexibility causally relates to other traits will allow researchers to develop
732 robust theory about the mechanisms and functional impact of flexibility, and when to invoke it as a primary
733 driver in a given context, such as a rapid geographic range expansion. Indeed, we are already in the process
734 of testing the latter hypothesis by conducting cross-population research on great-tailed grackles to test
735 whether a population on the range edge is more flexible (Logan CJ et al., 2020). That we were able to
736 manipulate flexibility, which had causal effects on flexible behavior in a different context (multi-access box)
737 as well as a different cognitive ability (innovativeness), demonstrates that flexibility manipulations could
738 be useful in training individuals of other species in how to be more flexible. This could have important
739 implications for threatened and endangered taxa (such as informing the choice of individuals for captive
740 breeding or introduction programs where individuals or their offspring are released into novel areas), as well
741 as for habituating zoo animals or other managed populations to novelty. If such a flexibility manipulation
742 was successful, it could then change their behavior in this and other domains, giving them a better chance of
743 succeeding in human modified environments. This is the focus of our new research program, ManyIndividuals,
744 where we manipulate flexibility using serial reversals in the wild in species that are successful and at risk
745 and determine whether the manipulation improves their success in human modified environments (Logan et
746 al., 2022).

747 ETHICS

748 This research is carried out in accordance with permits from the:

- 749 1) US Fish and Wildlife Service (scientific collecting permit number MB76700A-0,1,2)
- 750 2) US Geological Survey Bird Banding Laboratory (federal bird banding permit number 23872)
- 751 3) Arizona Game and Fish Department (scientific collecting license number SP594338 [2017], SP606267
752 [2018], and SP639866 [2019])
- 753 4) Institutional Animal Care and Use Committee at Arizona State University (protocol number 17-1594R)
- 754 5) University of Cambridge ethical review process (non-regulated use of animals in scientific procedures:
755 zoo4/17 [2017])

756 AUTHOR CONTRIBUTIONS

757 **Logan:** Hypothesis development, protocol development, data collection, data analysis and interpretation,
758 write up, revising/editing, materials/funding.

759 **Lukas:** Hypothesis development, simulation development, data interpretation, revising/editing.

760 **Blaisdell:** Prediction revision, assisted with programming the reversal learning touchscreen experiment,
761 protocol development, data interpretation, revising/editing.

762 **Johnson-Ulrich:** Prediction revision, programming, data collection, data interpretation, revising/editing.

763 **MacPherson:** Data collection, data interpretation, revising/editing.

764 **Seitz:** Prediction revision, programmed the reversal learning touchscreen experiment, protocol development,
765 data interpretation, revising/editing.

766 **Sevchik:** Data collection, revising/editing.

767 **McCune:** Added MAB log experiment, protocol development, data collection, data interpretation, revis-
768 ing/editing, materials.

769 **FUNDING**

770 This research is funded by the Department of Human Behavior, Ecology and Culture at the Max Planck Insti-
771 tute for Evolutionary Anthropology (2017-current), and by a Leverhulme Early Career Research Fellowship
772 to Logan (2017-2018).

773 **CONFLICT OF INTEREST DISCLOSURE**

774 We, the authors, declare that we have no financial conflicts of interest with the content of this article. CJ
775 Logan and D Lukas are Recommenders at PCI Ecology, and Logan used to be on the Managing Board
776 (2018-2022).

777 **ACKNOWLEDGEMENTS**

778 We thank our PCI Ecology recommender, Aurelie Coulon, and reviewers, Maxime Dahirel and Andrea
779 Griffin, for their feedback on the preregistration and post-study manuscript; Kevin Langergraber for serving
780 as our ASU IACUC PI; Ben Trumble and Angela Bond for logistical support; Melissa Wilson for sponsoring
781 our affiliations at Arizona State University and lending lab equipment; Kristine Johnson for technical advice
782 on great-tailed grackles; Arizona State University School of Life Sciences Department Animal Care and
783 Technologies for providing space for our aviaries and for their excellent support of our daily activities;
784 Julia Cissewski for tirelessly solving problems involving financial transactions and contracts; Sophie Kaube
785 for logistical support; Richard McElreath for project support; Aaron Blackwell and Ken Kosik for being the
786 UCSB sponsors of the Cooperation Agreement with the Max Planck Institute for Evolutionary Anthropology;
787 Tiana Lam, Anja Becker, and Brynna Hood for interobserver reliability video coding; Sawyer Lung for field
788 support; Alexis Breen for coding multi-access box videos; and our research assistants: Aelin Mayer, Nancy
789 Rodriguez, Brianna Thomas, Aldora Messinger, Elysia Mamola, Michael Guillen, Rita Barakat, Adriana
790 Boderash, Olateju Ojekunle, August Sevchik, Justin Huynh, Jennifer Berens, Amanda Overholt, Michael
791 Pickett, Sam Munoz, Sam Bowser, Emily Blackwell, Kaylee Delcid, Sofija Savic, Brynna Hood, Sierra
792 Planck, and Elise Lange.

793 SUPPLEMENTARY MATERIAL 1: Ability to detect actual effects

794 To begin to understand what kinds of effect sizes we will be able to detect given our sample size limitations
795 and our interest in decreasing noise by attempting to measure it, which increases the number of explanatory
796 variables, we used G*Power (v.3.1, Faul et al., 2007, 2009) to conduct power analyses based on confidence
797 intervals. G*Power uses pre-set drop down menus and we chose the options that were as close to our
798 analysis methods as possible (listed in each analysis below). Note that there were no explicit options for
799 GLMs (though the chosen test in G*Power appears to align with GLMs) or GLMMs or for the inclusion of
800 the number of trials per bird (which are generally large in our investigation), thus the power analyses are only
801 an approximation of the kinds of effect sizes we can detect. We realize that these power analyses are not fully
802 aligned with our study design and that these kinds of analyses are not appropriate for Bayesian statistics
803 (e.g., our MCMCglmm below), however we were unaware of better options at that time. Additionally, it
804 is difficult to run power analyses because it is unclear what kinds of effect sizes we should expect due to the
805 lack of data on this species for these experiments.

806 To address the power analysis issues, we ran simulations on our Arizona data set before conducting any
807 analyses in this preregistration.

808 **Planned:** We will first run null models (i.e., dependent variable $\sim 1 +$ random effects), which will allow us
809 to determine what a weak versus a strong effect is for each model. Then we will run simulations based on
810 the null model to explore the boundaries of influences (e.g., sample size) on our ability to detect effects of
811 interest of varying strengths. If simulation results indicate that our Arizona sample size is not larger than
812 the lower boundary, we will continue these experiments at the next field site until we meet the minimum
813 suggested sample size.

- 814 • **Implementation of the plan:** simulations were conducted in April 2020 (pre-data analysis) following
815 procedures in McElreath (2018). This meant that there were no null models because the simulations
816 using the full models are used to determine whether one can detect differences between effect sizes.

817 To run the simulations, we first constructed a **hypothesis-appropriate mathematical model** that en-
818 compasses the relationship between the variables of interest for each analysis: 1) number of loci solved on
819 the multi-access box \sim trials to reverse, and 2) latency to attempt a new locus on the multi-access box \sim
820 trials to reverse.

821 **Simulation and model: number of loci solved on the multi-access box \sim trials to reverse**

822 The model takes the form of:

$$823 \text{locisolved} \sim \text{Binomial}(4, p) \text{ [likelihood]}$$

$$824 \text{logit}(p) \sim \alpha[\text{batch}] + \beta\text{trials} \text{ [model]}$$

825 locisolved is the number of loci solved on the multi-access box, 4 is the total number of loci on the multi-
826 access box, p is the probability of solving any one locus across the whole experiment, α is the intercept and
827 each batch gets its own, β is the expected amount of change in locisolved for every one unit change in trials,
828 and trials is the number of trials to reverse a color preference.

829 Expected values for the number of loci solved on the multi-access box were set to either 2 or 0 (out of
830 4 loci maximum) because we were unsure of whether the grackles would be able to solve any loci on the
831 multi-access box because this experiment had never been done on this species before. Expected values for
832 reversal learning using colored tubes (mean, standard deviation, and range of number of trials to reverse a
833 color preference) were based on previously published data on great-tailed grackles (Logan, 2016). This data
834 indicates that the average number of trials to reverse a preference is 91 and the standard deviation is 21. In
835 our model, the variation in the actual data is reflected by both the population standard deviation and the
836 expected amount of change related to the explanatory variable. After running simulations, we identified the
837 following distributions and priors to be the most likely for our expected data:

$$838 \alpha \sim \text{Normal}(4,10) \text{ [}\alpha \text{ prior]}$$

839 $\beta \sim \text{Normal}(0,5)$ [β prior]

840 We used normal distributions for α and β because they are (or are based on) sums with large means (see
841 Figure 10.6 in McElreath, 2018). For the β prior, we had no expectation about whether the relationship
842 would be positive or negative, therefore we centered it on 0 (the mean).

843 **Simulation and model: latency to attempt a new locus on the multi-access box ~ trials to**
844 **reverse**

845 For the average latency to attempt a new locus on the multi-access box as it relates to trials to reverse (both
846 are measures of flexibility), we simulated data and set the model as follows:

847 latency \sim gamma-Poisson(λ_i, ϕ) [*likelihood*]

848 $\log(\lambda_i) \sim \alpha[\text{batch}] + \beta\text{trials}$ [*the model*]

849 latency is the average latency to attempt a new locus on the multi-access box, λ_i is the rate (probability of
850 attempting a locus in each second) per bird (and we take the log of it to make sure it is always positive; birds
851 with a higher rate have a smaller latency), ϕ is the dispersion of the rates across birds, α is the intercept
852 for the rate per batch, β is the expected amount of change in the rate of attempting to solve in any given
853 second for every one unit change in trials, and trials is the number of trials to reverse a color preference.

854 Expected values for the latency to attempt a new locus on the multi-access box was set to between 1-2700
855 sec because the experiment ends for a bird if they do not obtain the food in 3 consecutive trials, and each
856 trial can last up to 15 min. Because we did not have prior data for this species on this test, we set the mean
857 to 300 sec, which is half way through a usual 10 min trial because it seems likely that if a bird is going to
858 attempt another locus, it will likely do so at the next opportunity, especially after being successful in the
859 previous trial. Expected values for reversal learning using colored tubes are the same as above. After running
860 simulations, we identified the following to be the most likely distributions and priors for our expected data:

861 $\phi \sim 1/(\text{Exponential}(1))$ [ϕ prior]

862 $\alpha \sim \text{Normal}(300,50)$ [α prior]

863 $\beta \sim \text{Normal}(0,5)$ [β prior]

864 We used a gamma-Poisson distribution for latency because it constrains the values to be positive and to
865 primarily occur sooner rather than later, which is what we expect from the grackles (based on data from New
866 Caledonian crows and kea in Auersperg et al., 2011). For ϕ , we used an exponential distribution because it
867 is standard for this paramter. We used normal distributions for α and β because they are (or are based on)
868 sums with large means (see Figure 10.6 in McElreath, 2018). For the β prior, we had no expectation about
869 whether the relationship would be positive or negative, therefore we centered it on 0 (the mean).

870 We translated the simulation output into effect sizes and examined what kind of effect size these
871 parameter values represent (Table SM1.1). For each β , we calculated the effect size (Box 13.3 in Lajeunesse
872 et al., 2013: linear regression):

873
$$r = \beta (\text{SDx} / \text{SDy}) = \beta (1.5 / 21)$$

874 Where r is the Pearson product moment correlation and SD is the standard deviation. For the standard
875 deviation of x (number of loci solved on the multiaccess box), we estimated a possible value of 1.5. For the
876 standard deviation of y (trials to reverse), we used 21 from the Santa Barbara grackle data (Logan, 2016).
877 We then calculated the effect sizes and R^2 values for each value of β .

	Beta	SDx	SDy	Effect size	\hat{R}^2
878	-5	1.5	21	-0.357	0.128
	-1	1.5	21	-0.071	0.005
879	0	1.5	21	0.000	0.000

880 **Table SM1.1.** The connection between β and effect sizes (SDx =standard deviation of x , which is the number
881 of loci solved; SDy =standard deviation of y , which is the number of trials to reverse; R^2 = R squared).

882 We then used the simulations to run **models** on simulated data to estimate the measurement error associated
883 with varying sample size, β , and the range of multi-access box loci solved or latency to attempt a new locus
884 (Table SM1.2). Before running the models, we decided that a model would detect an effect if 89% of the
885 posterior sample was on the same side of zero (following McElreath, 2018). We ran the simulation with
886 $\beta=3$ (latency) because this was a high value at which an appropriate range of values were observed in the
887 simulation testing phase, $\beta=0$ because this would be the scenario in which there is no relationship between
888 the response variable and the trials to reverse, and $\beta=-1$ to determine how small of a difference we can detect
889 and with what amount of associated noise (σ). Sigma (σ) is the standard deviation in the trials to reverse
890 if the trials to reverse is a normal distribution. In all simulations, the mean in the trials to reverse was set
891 to 91. Therefore, a (σ) of 14 is 15% noise (14/91). We found that when (σ) is larger than 14, we cannot
892 detect even the largest effect of trials to reverse on loci solved or latency because there are some simulations
893 where the estimated regression coefficient crosses zero. When $\beta=0$ we want all of the regression coefficients
894 to cross zero (10 out of 10 random repetitions) and when $\beta \neq 0$ we want none of the regression coefficients
895 to cross zero (0 out of 10 random repetitions). We ran the models several times with various parameters to
896 determine at what point this was the case for each combination of parameters.

Beta	n	Sigma	Regression coefficient crosses zero	Regression coefficient	Range of MAB loci solved
-5	15	15	1/10	-5.90	0-4
-5	15	14	0/10	-5.11	0-4
-5	15	12	0/10	-4.79	0-4
-5	15	10	0/10	-4.31	0-4
-5	10	10	1/10	-4.35	0-4
-5	10	9	0/10	-5.26	0-4
-5	8	10	1/10	-5.35	0-4
-5	8	9	0/10	-4.22	0-4
-5	8	8	0/10	-3.08	0-4
-5	8	8	1/10	-4.74	0-2
-5	8	7	3/10	-6.74	0-2
-5	8	5	0/10	-3.08	0-2
-5	10	9	3/10	-4.51	0-2
-5	10	7	1/10	-7.67	0-2
-5	10	6	2/10	-5.16	0-2
-5	10	5	1/10	-4.57	0-2
-5	10	4	0/10	-5.02	0-2
-5	15	14	2/10	-3.07	0-2
-5	15	13	5/10	1.68	0-2
-5	15	10	5/10	-8.20	0-2
-5	15	8	3/10	-4.01	0-2
-5	15	6	0/10	-6.03	0-2
-5	15	7	1/10	-8.06	0-2
0	15	14	10/10	-3.23	0-2
0	15	14	10/10	0.43	0-4
-1	15	14	10/10	-1.53	0-4
-1	15	10	10/10	-0.73	0-4
-1	15	5	3/10	0.19	0-4
-1	15	3	1/10	0.18	0-4
-1	15	2	0/10	-1.07	0-4
-1	15	2	3/10	-1.67	0-2
-1	15	1	1/10	-1.12	0-2

899 **Table SM1.2.** Simulation outputs from varying β , sample size (n), σ , and whether the actual range of
900 multi-access box [MAB] loci solved were 0-2 or 0-4 (we did not know how many loci the grackles would
901 be able to solve before we started collecting data so we ran two simulations. The grackles ended up being
902 able to solve all four loci on both multi-access boxes, therefore we must use only those rows associated with
903 “Range of MAB loci solved” = 0-4). This table is useful for the analyses involving the number of loci solved
904 on the multi-access box, but not the latency to switch to attempting a new locus on the multi-access box,
905 which uses a different (gamma poisson) model.

906 This shows that we would have the power to detect a medium effect (-0.357 in Table M1) with a sample
907 size of 15 if the noise (σ) is <15%. We would be unlikely to get a false negative because there were no false
908 negatives in the simulations (i.e., the posterior sample range did not cross zero). With this sample size, when
909 $\beta=0$, there are no false positives (i.e., the posterior sample range always included zero). However, we would
910 not be able to detect a weak effect unless the noise (σ) was much smaller.

911 **SUPPLEMENTARY MATERIAL 2: Interobserver reliability of dependent vari-**
912 **ables (unregistered analyses)**

913 To determine whether experimenters coded the dependent variables in a repeatable way, hypothesis-blind
914 video coders were first trained in video coding the dependent variable, and then they coded at least 20% of
915 the videos in the reversal (tubes) and multi-access box experiments. We randomly chose a subset of all of
916 the birds who participated in each experiment using random.org:

- 917 • Reversal 6/20 grackles (30% with half from the control group): Chalupa, Avocada, Diablo, Fideo,
918 Tomatillo, Adobo
- 919 • Multi-access box plastic 3/15 grackles (20%): Habanero, Queso, Chalupa
- 920 • Multi-access box log 3/12 grackles (25%): Diablo, Adobo, Yuca

921 Video coders then analyzed all videos from these birds. The experimenter's data was compared with the
922 video coder data using the intra-class correlation coefficient (ICC) to determine the degree of bias in the
923 regression slope (Hutcheon et al. (2010), using the irr package in R: Gamer et al. (2012)). Note that the
924 data in columns from coders 1 and 2 in the data sheets were aligned based on similar numbers between
925 coders to prevent disagreements near the top of the data sheet from misaligning all subsequent entries.

926 **Interobserver reliability training** To pass **interobserver reliability (IOR) training**, video coders
927 needed an ICC score of 0.90 or greater to ensure the instructions were clear and that there was a high degree
928 of agreement across coders (see R code comments for details).

929 **Alexis Breen** (compared with experimenter's live coding):

- 930 • Multi-access box: correct choice unweighted Cohen's Kappa=0.90 (confidence boundaries=0.77-1.00,
931 n=33 data points)
- 932 • Multi-access box: locus solved unweighted Cohen's Kappa=0.90 (confidence boundaries=0.76-1.00,
933 n=33 data points)

934 Note: Breen was not a hypothesis-blind video coder. She contributed to extensive video coding across
935 the whole project, however, for interobserver reliability analyses, her data were always compared with a
936 hypothesis-blind coder's data.

937 **Anja Becker** (compared with experimenter's live coding):

- 938 • Reversal: correct choice ICC=1.00 (confidence boundaries=1.00-1.00, n=25 data points)

939 **Tiana Lam** (compared with experimenter's live coding):

- 940 • Multi-access box: correct choice ICC=0.90 (confidence boundaries=0.77-1.00, n=33 data points)
- 941 • Multi-access box: locus solved unweighted Cohen's Kappa=0.95 (confidence boundaries=0.84-1.00,
942 n=33 data points)

943 **Brynna Hood** (compared with experimenter's live coding):

- 944 • Multi-access log: correct choice unweighted Cohen's Kappa=1.00 (confidence boundaries=1.00-1.00,
945 n=29 data points)
- 946 • Multi-access log: locus solved unweighted Cohen's Kappa=1.00 (confidence boundaries=1.00-1.00,
947 n=29 data points)

948 **Interobserver reliability** Interobserver reliability scores (minimum 20% of the videos) were as follows:

949 *Brynna Hood* (compared with experimenter's live coding):

- 950 • Multi-access log: correct choice unweighted Cohen's Kappa=0.91 (confidence boundaries=0.76-1.00,
951 n=39 data points)
- 952 • Multi-access log: locus solved unweighted Cohen's Kappa=1.0 (confidence boundaries=1.0-1.00, n=39
953 data points)

954 *Tiana Lam* (compared with experimenter's live coding):

- 955 • Multi-access box: correct choice unweighted Cohen's Kappa=0.83 (confidence boundaries=0.73-0.92,
956 n=102 data points)
- 957 • Multi-access box: locus solved unweighted Cohen's Kappa=0.90 (confidence boundaries=0.830-0.97,
958 n=102 data points)

959 *Anja Becker* (compared with experimenter's live coding):

- 960 • Reversal: correct choice ICC=0.99 (confidence boundaries=0.98-0.99, n=3280 data points)

961 These scores indicate that the dependent variables are repeatable to a high or extremely high degree given
962 our instructions and training

963 **SUPPLEMENTARY MATERIAL 3: Prediction 2 model outputs**

964 **Table SM3.** Model outputs for the number of loci solved and the latency to switch loci after passing
965 criterion on a different locus on the plastic (models 1-5 and 9-11) and wooden (models 6-8 and 12-14) multi-
966 access boxes. SD=standard deviation, the 89% prediction intervals are shown, n_eff=effective sample size,
967 Rhat4=an indicator of model convergence (1.00 is ideal), b=the slope of the relationship between loci solved
968 or average switch latency and the number of trials to pass the reversal.

	Mean	SD	Lower 89 percentile compatibility interval (5.5%)	Upper 89 percentile compatibility interval (94.5%)	n_eff	Rhat4
MODEL 1 (last reversal): loci solved plastic - a[batch] + b*trials						
a[1]	0.04	0.46	-0.70	0.78	2304	1.00
a[2]	0.29	0.36	-0.30	0.87	2456	1.00
a[3]	-0.78	0.55	-1.65	0.08	2510	1.00
b	-0.22	0.25	-0.63	0.18	2364	1.00
MODEL 2 (last reversal): loci solved plastic - a + b*trials						
a	-0.02	0.24	-0.40	0.35	1466	1.00
b	-0.46	0.31	-0.97	-0.01	1383	1.00
MODEL 3 (last reversal): trials - a[batch]						
a[1]	0.09	0.37	-0.48	0.69	2095	1.00
a[2]	-0.21	0.29	-0.68	0.25	1715	1.00
a[3]	0.25	0.39	-0.38	0.86	2161	1.00
sigma	1.03	0.21	0.75	1.39	2049	1.00
MODEL 4: loci solved - a[condition]						
a[1] control	-0.11	0.32	-0.62	0.40	1311	1.00
a[2] manipulated	0.15	0.39	-0.46	0.80	1222	1.00
MODEL 5 (first reversal): loci solved plastic - a + b*trials						
a	0.00	0.24	-0.37	0.39	1208	1.00
b	-0.44	0.30	-0.94	0.02	1273	1.00
MODEL 6 (last reversal): loci solved wooden - a + b*trials						
a	1.06	0.27	0.63	1.50	1255	1.00
b	0.41	0.43	-0.21	1.13	1107	1.00
MODEL 7: loci solved - a[condition]						
a[1] control	-0.45	0.40	-1.10	0.18	1161	1.00
a[2] manipulated	0.77	0.41	0.13	1.44	1302	1.00
MODEL 8 (first reversal): loci solved wooden - a + b*trials						
a	0.11	0.26	-0.30	0.52	1221	1.00
b	-0.50	0.35	-1.09	0.04	1234	1.00
MODEL 9 (last reversal): avg switch latency plastic - a + b*trials						
a	4.93	0.30	4.45	5.41	1235	1.01
b	0.46	0.29	0.00	0.92	1363	1.00
phi	0.93	0.35	0.44	1.55	1476	1.00
MODEL 10: avg switch latency plastic - a[condition]						
a[1] manipulated	4.07	0.39	3.46	4.68	1027	1.00
a[2] control	5.18	0.39	4.50	5.76	1006	1.00
phi	0.91	0.41	0.37	1.63	925	1.01
MODEL 11 (first reversal): avg switch latency plastic - a + b*trials						
a	4.93	0.29	4.46	5.39	1488	1.00
b	0.46	0.28	0.02	0.93	1211	1.00
phi	0.94	0.36	0.44	1.60	1447	1.00
MODEL 12 (last reversal): avg switch latency wooden - a + b*trials						
a	5.75	0.28	5.28	6.18	1049	1.00
b	-0.41	0.32	-0.86	0.15	1281	1.01
phi	1.04	0.42	0.48	1.77	1456	1.00
MODEL 13: avg switch latency wooden - a[condition]						
a[1] control	5.31	0.42	4.61	5.95	701	1.00
a[2] manipulated	5.34	0.44	4.61	6.00	620	1.01
phi	0.66	0.32	0.25	1.25	806	1.00
MODEL 14 (first reversal): avg switch latency wooden - a + b*trials						
a	5.71	0.26	5.28	6.12	1109	1.00
b	-0.50	0.28	-0.89	-0.01	1308	1.00
phi	1.08	0.41	0.53	1.80	1347	1.00

969

970

971 **SUPPLEMENTARY MATERIAL 4: Reversal learning experiments: discrimi-**
972 **nating shapes on the touchscreen compared with color using tubes**

973 In the tube experiment, it took four grackles an average of 40 trials (sd=12) in the initial discrimination
974 phase to learn to prefer a color, while it took the same individuals an average of 390 trials (sd=59) to learn
975 to prefer a shape using the touchscreen (Queso, Mole, Habanero, and Tapa). The two individuals who were
976 faster to learn in the tube experiment were slower to learn in the touchscreen experiment. For the reversal,
977 it took three of these individuals (Queso, Mole, and Habanero) an average of 80 trials (sd=14) to reverse
978 their colored tube preference, and an average of 362 trials (sd=111) to reverse their shape preference on the
979 touchscreen (Tapa had to be released back to the wild before finishing the experiment, but was on trial 629
980 in reversal one of the touchscreen experiment at the time of release. In the tube experiment, she was also
981 the slowest of the four to reverse at 100 trials). All three individuals were about equally fast at the reversal
982 in the tube experiment, while their reversal learning speeds differed on the touchscreen. The touchscreen
983 training data and a summary of the training process is detailed in Seitz et al. (2021).

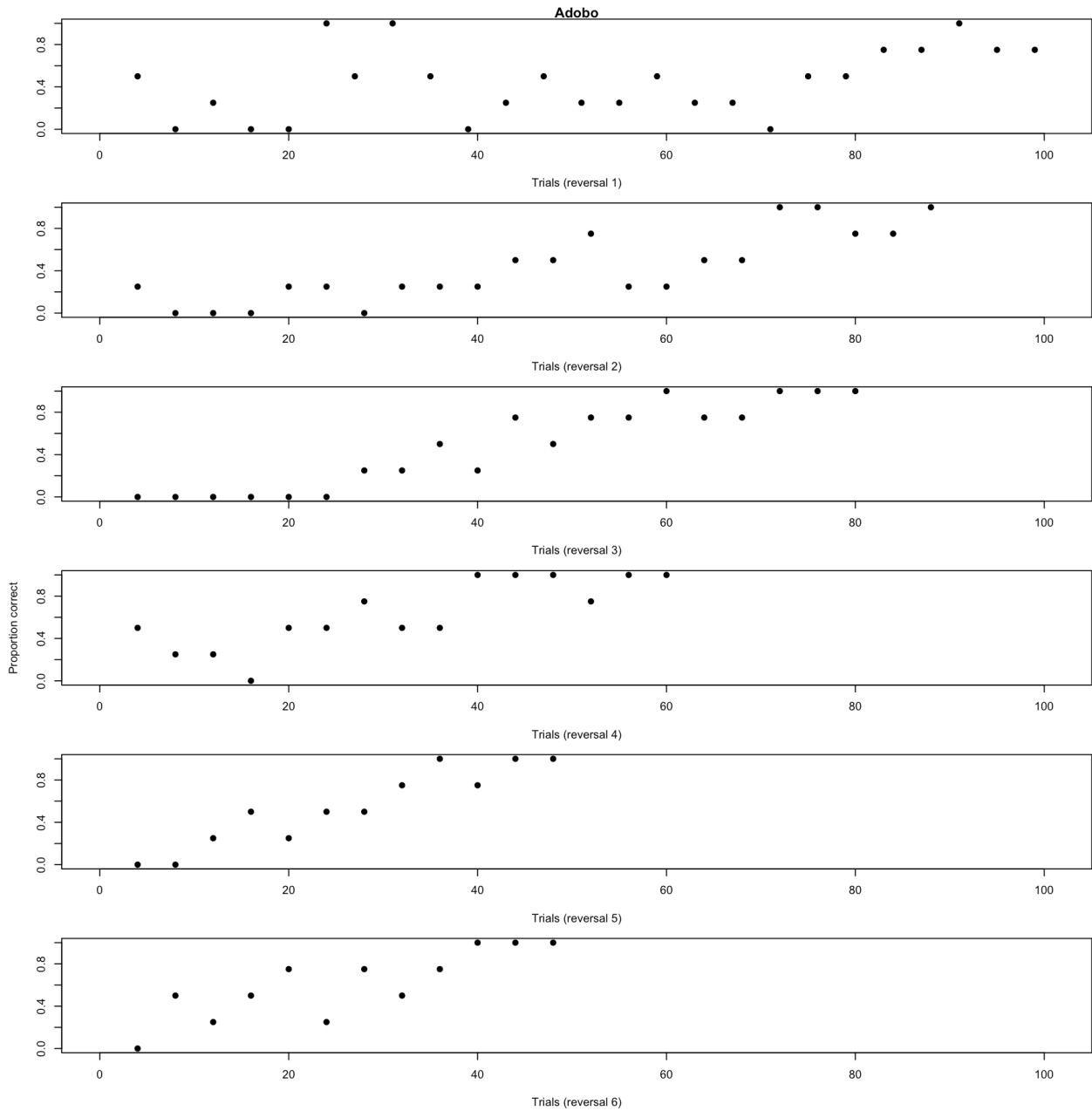
⁹⁸⁴ **SUPPLEMENTARY MATERIAL 5: Summarized results per bird**

985 **Table SM5.** Summarized results per bird in the reversal learning (tube and touchscreen) and multi-access box (plastic and wooden) experiments.
986 “Reversals to pass” indicates how many serial reversals it took a bird to pass criterion (passing two consecutive reversals in 50 trials or less) if they
987 were in the flexibility manipulation condition. X indicates the bird attempted, but did not pass that experiment. Note: Tapa did not finish the MAB
988 log experiment; Marisco’s MAB log experiment ended too early due to experimenter error (timed out on 2 consecutive sessions, not 3); Mole and
989 Habanero: do not count MAB plastic number of options solved because they were given the box fully put together for habituation due to experimenter
990 error; Taco was the first juvenile we tested and we did not put him in the flexibility experiment: he received 1 reversal and moved on to his next test,
991 therefore he was essentially a control bird without the matched yellow tube experience.

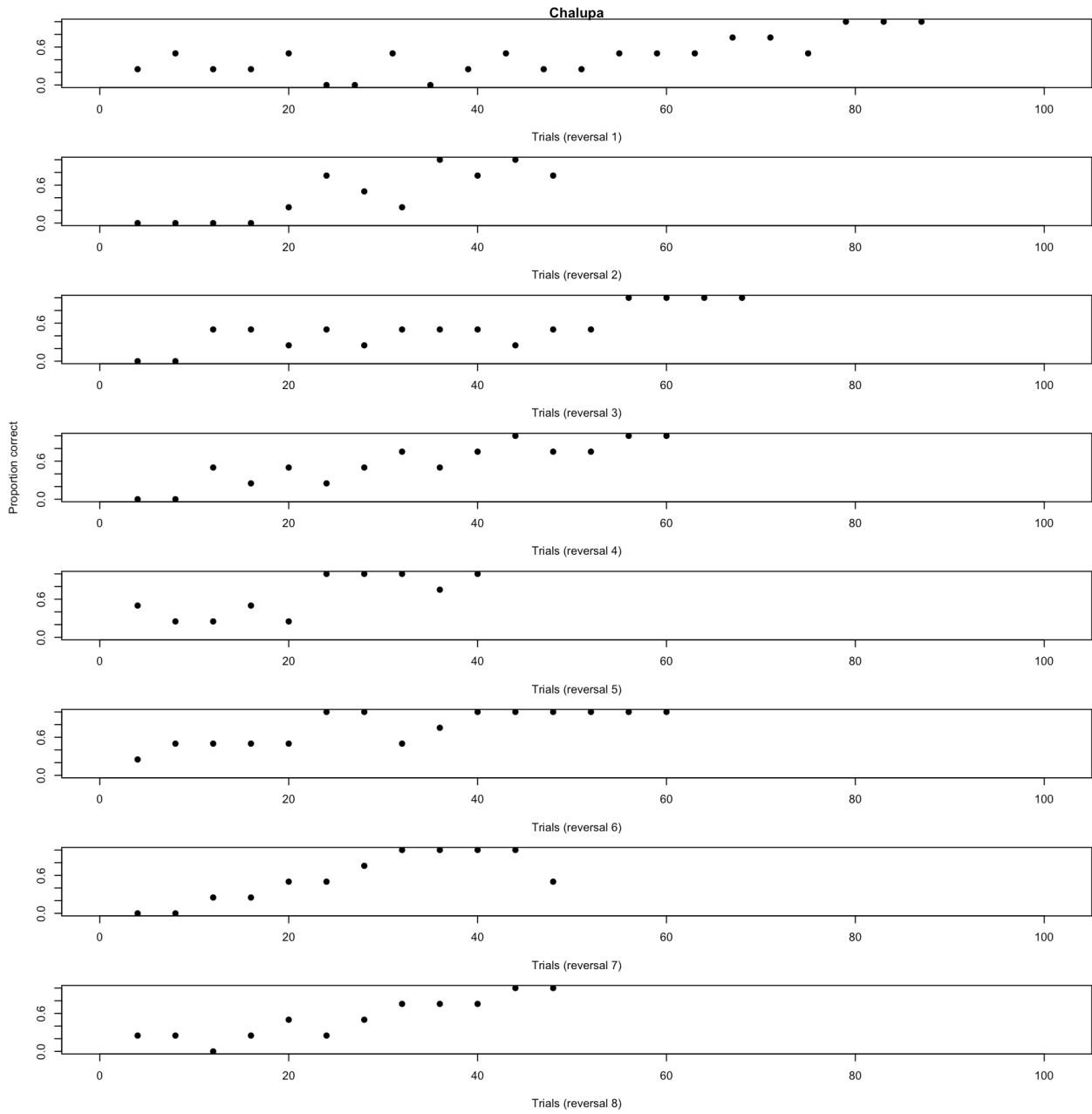
Bird	Batch	Sex	Trials to learn (tube)	Trials to first reversal (tube)	Trials to last reversal (tube)	Reversals to pass	Total loci solved (MAB plastic)	Total loci solved (MAB wooden)	Average latency to attempt new locus (MAB plastic)	Average latency to attempt new locus (MAB wooden)	Trials to learn (touchscreen)	Trials to first reversal (touchscreen)	Motor actions (MAB plastic)	Motor actions (MAB wooden)
Tomatillo	1	M	40	50	50	Control	3		317				13	
Queso	1	M	50	70	70	Control	1		88		330	460	8	
Tapa	1	F	30	100	100	Control	4		685		450	(629+)	12	
Yuca	3	F	40	80	80	Control	4	4	132	77			13	16
Marisco	3	M	40	50	50	Control	1	2		208			3	7
Pizza	3	M	50	60	60	Control	0	1		1482			0	8
Mofongo	4	M	20	40	40	Control	3	4	502	630			13	14
Taquito	4	M	90	160	160	Control	0	4		100			11	10
Chalupa	1	F	50	90	50	8	0						6	
Mole	1	M	30	70	50	7	4	4	356	1173	431	307	14	15
Habanero	1	M	50	80	40	6	4		28		350	290	15	
Diablo	3	M	20	80	40	8	2	1	25				10	2
Burrito	3	M	40	60	23	8	3	4	76	391			17	18
Adobo	3	M	50	100	50	6	4	4	31	79			16	18
Chilaquile	3	JM	30	40	30	6	4	4	44	170			19	11
Pollito	4	M	40	60	40	8	0	3		668			0	11
Taco	3a	JM	50	80	80	(Control)	1	4		117			2	19
Memela	1	F	50	60	80	X (11+)								
Fideo	2	M	60	70	70	Control								
Avocada	1	F	50	100	100	Control								
Huachinago	3	M	70			Control								
Guacamole	4	M	30											

994 **SUPPLEMENTARY MATERIAL 6: Prediction 4 learning strategy figures**

995 Below are figures for the proportion of trials correct by trial number and reversal for each bird.

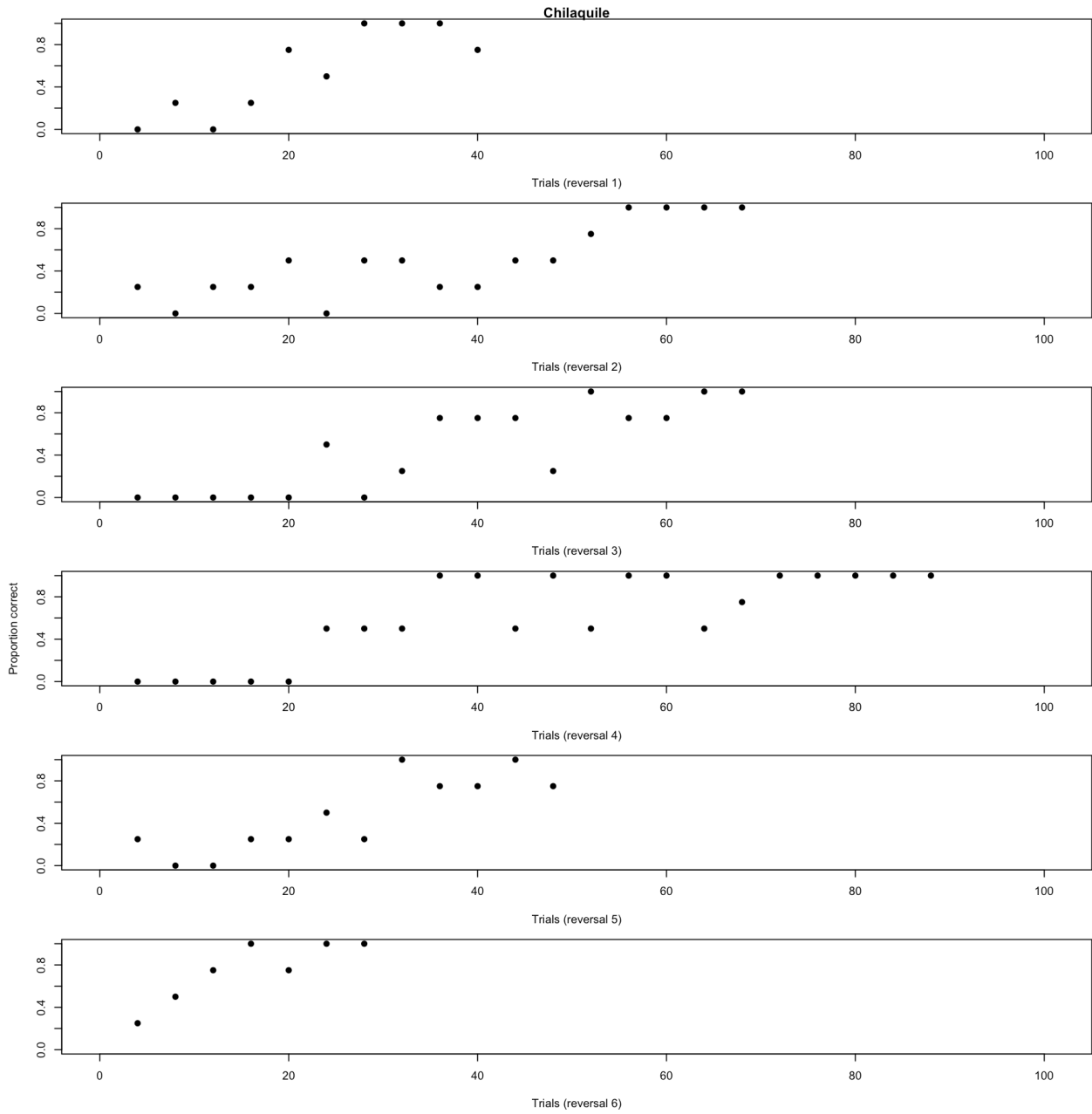


996
997 **Figure SM6.1.** Adobo's proportion of trials correct by trial number and reversal.



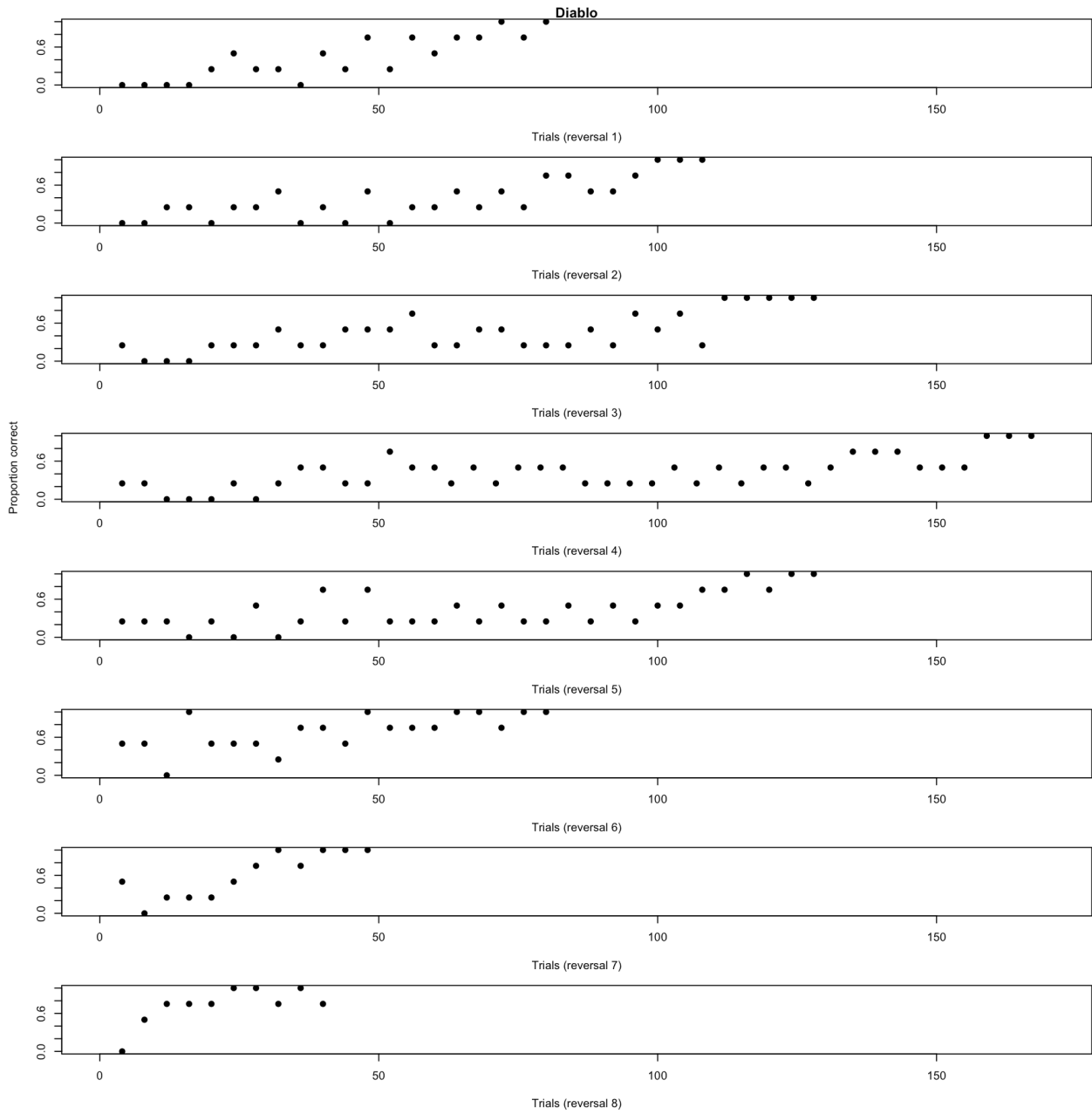
998

999 **Figure SM6.2.** Chalupa's proportion of trials correct by trial number and reversal.



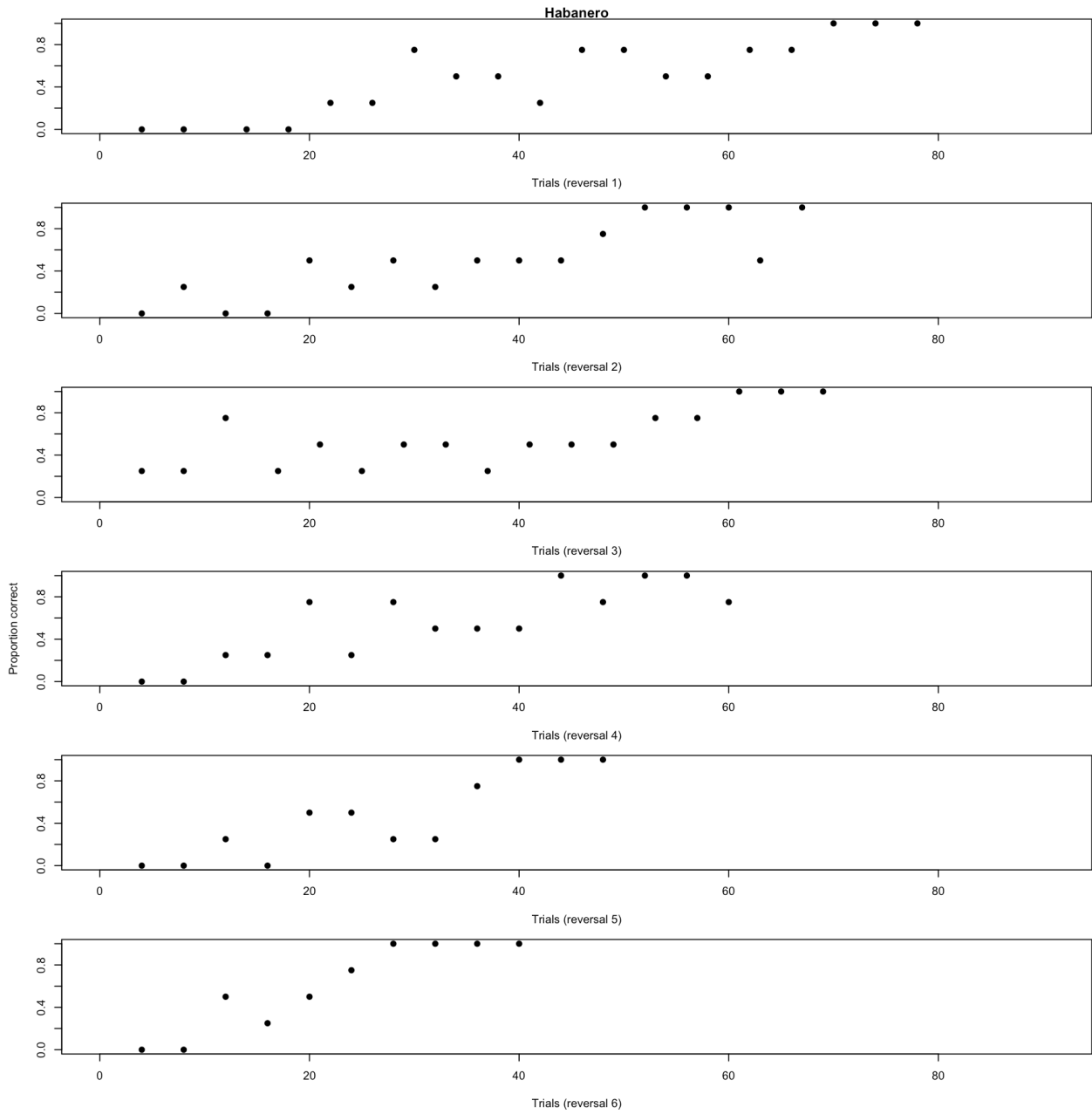
1000

1001 **Figure SM6.3.** Chilaquile's proportion of trials correct by trial number and reversal.



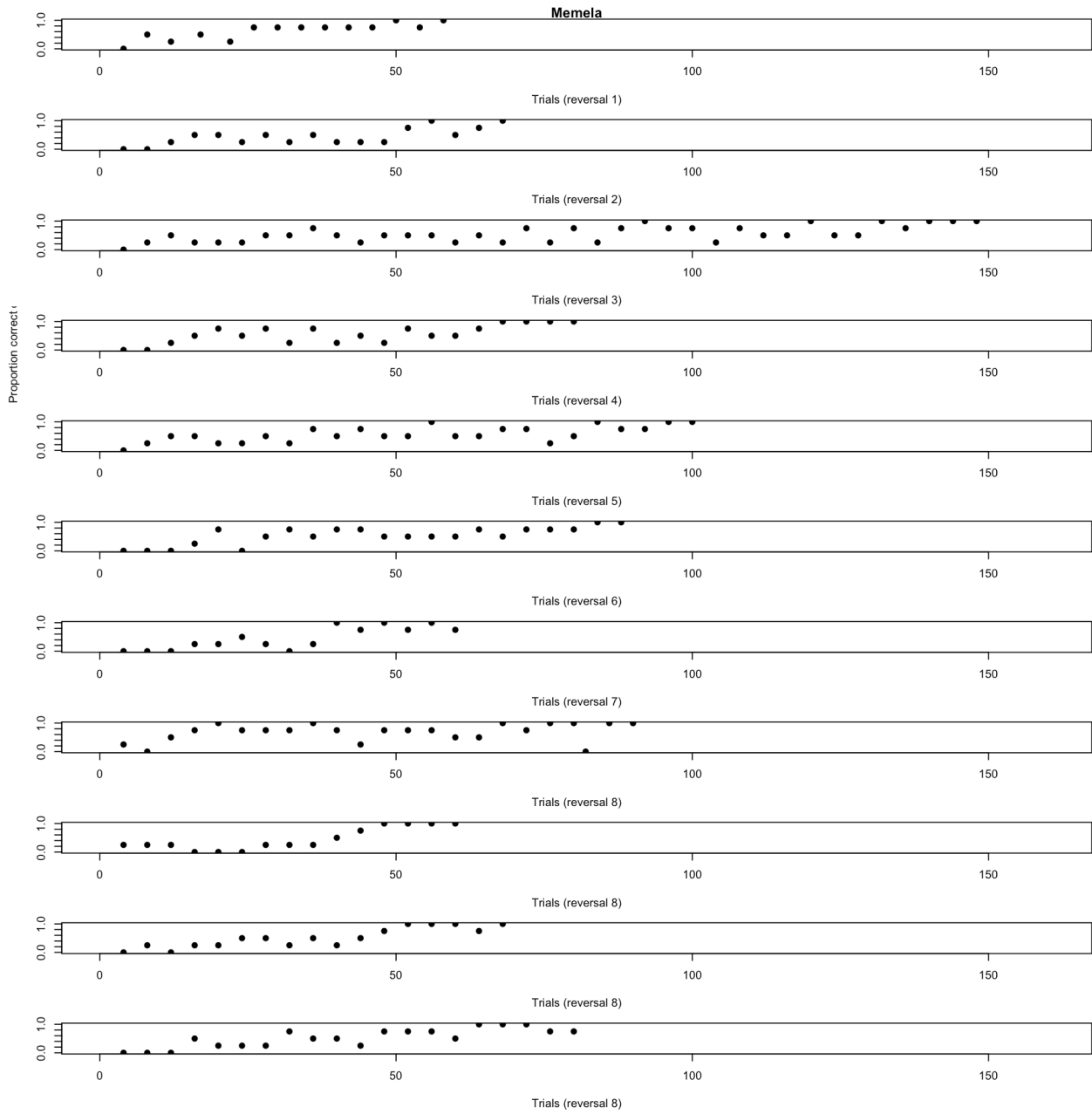
1002

1003 **Figure SM6.4.** Diablo's proportion of trials correct by trial number and reversal.



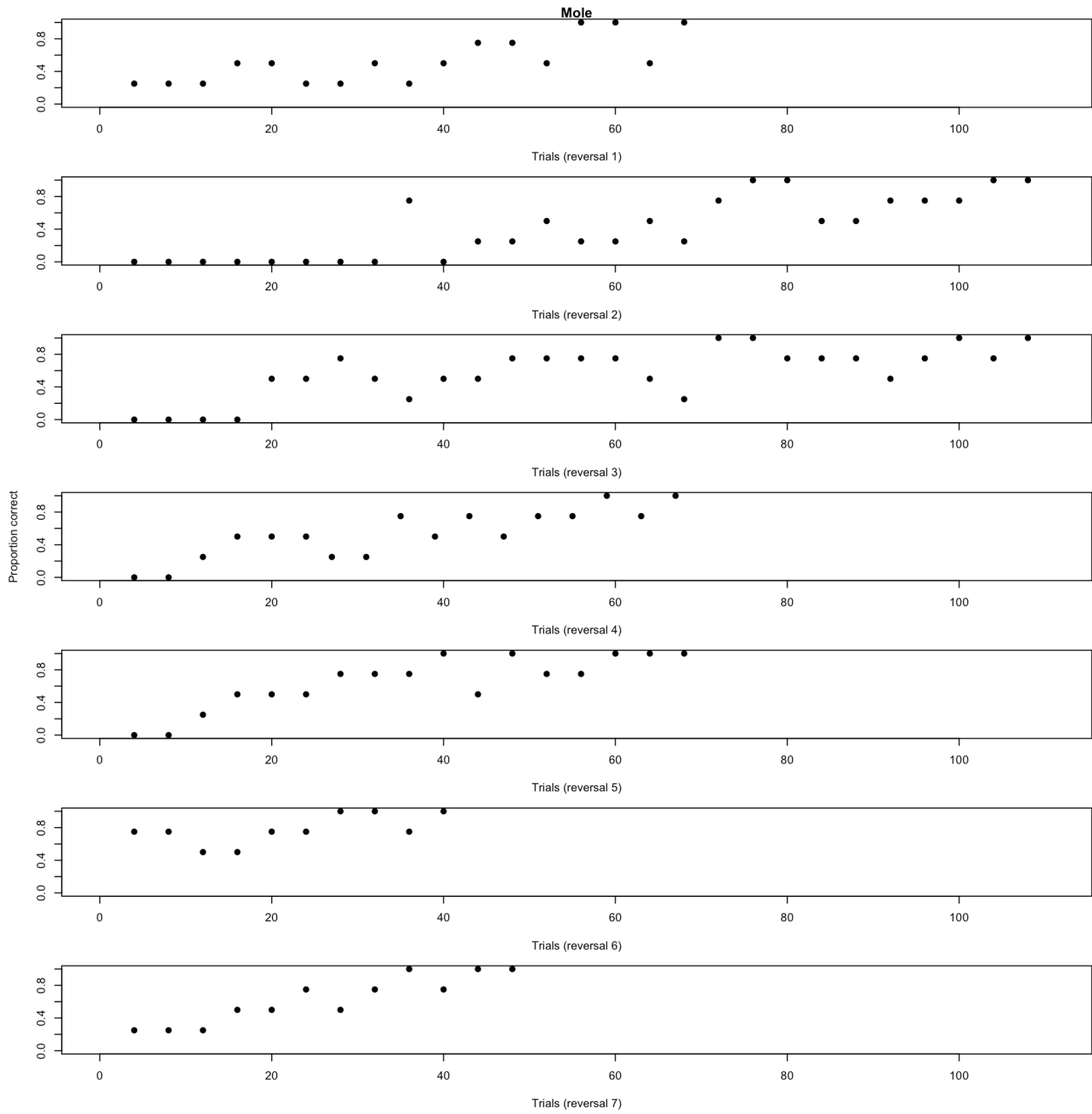
1004

1005 **Figure SM6.5.** Habanero's proportion of trials correct by trial number and reversal.



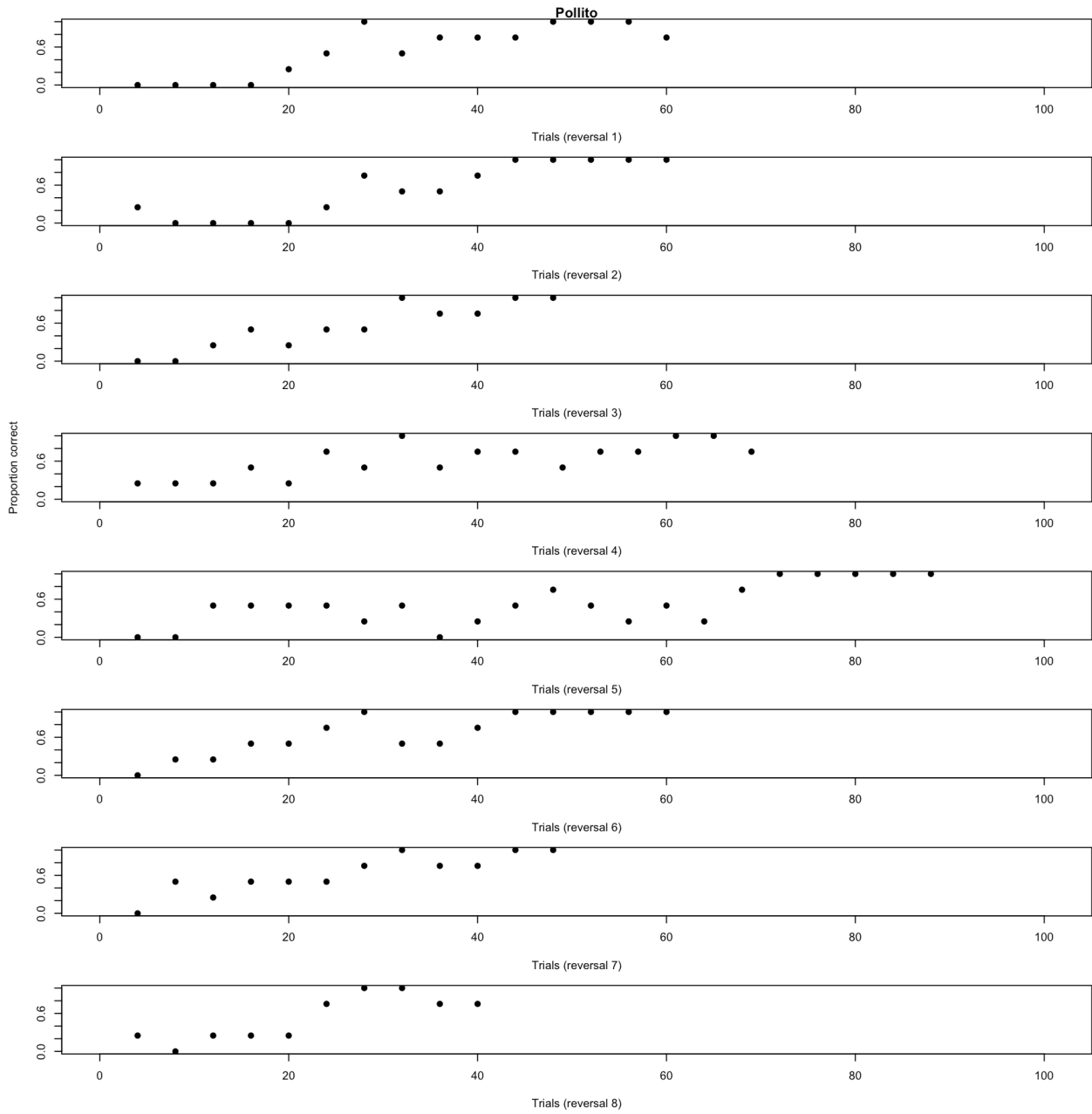
1006

1007 **Figure SM6.6.** Memela's proportion of trials correct by trial number and reversal.



1008

1009 **Figure SM6.7.** Mole's proportion of trials correct by trial number and reversal.



1010

1011 **Figure SM6.8.** Pollito's proportion of trials correct by trial number and reversal.

1012 REFERENCES

- 1013 Aplin, L. M., Farine, D. R., Morand-Ferron, J., Cockburn, A., Thornton, A., & Sheldon, B. C. (2015).
 1014 Experimentally induced innovations lead to persistent culture via conformity in wild birds. *Nature*,
 1015 *518*(7540), 538–541.
- 1016 Auersperg, A. M. I., Bayern, A. M. P. von, Gajdon, G. K., Huber, L., & Kacelnik, A. (2011). Flexibility in
 1017 problem solving and tool use of kea and New Caledonian crows in a multi access box paradigm. *PLOS*
 1018 *ONE*, *6*(6), e20231. <https://doi.org/10.1371/journal.pone.0020231>
- 1019 Bartoń, K. (2020). *MuMIn: Multi-model inference*. <https://CRAN.R-project.org/package=MuMIn>
- 1020 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4.
 1021 *Journal of Statistical Software*, *67*(1), 1–48. <https://doi.org/10.18637/jss.v067.i01>
- 1022 Bates, D., Maechler, M., & Bolker, B. (2012). *lme4: Linear mixed-effects models using S4 classes (2011)*. *R*
 1023 *package version 0.999375-42*.
- 1024 Bergstrom, C. T., & Lachmann, M. (2004). Shannon information and biological fitness. *Information Theory*
 1025 *Workshop, 2004. IEEE*, 50–54.
- 1026 Blaisdell, A. P., & Cook, R. G. (2005). Integration of spatial maps in pigeons. *Animal Cognition*, *8*(1), 7–16.
- 1027 Bürkner, P.-C., Gabry, J., Kay, M., & Vehtari, A. (2020). Posterior: Tools for working with posterior
 1028 distributions. *Earthquake Spectra, R Package Version 0.1*, *3*.
- 1029 Bussey, T. J., Padain, T. L., Skillings, E. A., Winters, B. D., Morton, A. J., & Saksida, L. M. (2008).
 1030 The touchscreen cognitive testing method for rodents: How to get the best out of your rat. *Learning &*
 1031 *Memory*, *15*(7), 516–523.
- 1032 Chow, P. K. Y., Lea, S. E., & Leaver, L. A. (2016). How practice makes perfect: The role of persistence,
 1033 flexibility and learning in problem-solving efficiency. *Animal Behaviour*, *112*, 273–283. [https://doi.org/](https://doi.org/10.1016/j.anbehav.2015.11.014)
 1034 [10.1016/j.anbehav.2015.11.014](https://doi.org/10.1016/j.anbehav.2015.11.014)
- 1035 Cook, R. G., Geller, A. I., Zhang, G.-R., & Gowda, R. (2004). Touchscreen-enhanced visual learning in rats.
 1036 *Behavior Research Methods, Instruments, & Computers*, *36*(1), 101–106.
- 1037 Diquelou, M. C., Griffin, A. S., & Sol, D. (2015). *The role of motor diversity in foraging innovations: A*
 1038 *cross-species comparison in urban birds*.
- 1039 Drayton, L. A., & Santos, L. R. (2014). Insights into intraspecies variation in primate prosocial behavior:
 1040 Capuchins (*cebus apella*) fail to show prosociality on a touchscreen task. *Behavioral Sciences*, *4*(2),
 1041 87–101.
- 1042 Eddelbuettel, D., & François, R. (2011). Rcpp: Seamless R and C++ integration. *Journal of Statistical*
 1043 *Software*, *40*(8), 1–18. <https://doi.org/10.18637/jss.v040.i08>
- 1044 Faul, F., Erdfelder, E., Buchner, A., & Lang, A.-G. (2009). Statistical power analyses using g* power
 1045 3.1: Tests for correlation and regression analyses. *Behavior Research Methods*, *41*(4), 1149–1160. [https://doi.org/](https://doi.org/10.3758/BRM.41.4.1149)
 1046 [10.3758/BRM.41.4.1149](https://doi.org/10.3758/BRM.41.4.1149)
- 1047 Faul, F., Erdfelder, E., Lang, A.-G., & Buchner, A. (2007). G* power 3: A flexible statistical power analysis
 1048 program for the social, behavioral, and biomedical sciences. *Behavior Research Methods*, *39*(2), 175–191.
 1049 <https://doi.org/10.3758/BF03193146>
- 1050 Federspiel, I. G., Garland, A., Guez, D., Bugnyar, T., Healy, S. D., Güntürkün, O., & Griffin, A. S. (2017).
 1051 Adjusting foraging strategies: A comparison of rural and urban common mynas (*acridotheres tristis*).
 1052 *Animal Cognition*, *20*(1), 65–74.
- 1053 Gabry, J., & Češnovar, R. (2021). *Cmdstanr: R interface to 'CmdStan'*.
- 1054 Gabry, J., Simpson, D., Vehtari, A., Betancourt, M., & Gelman, A. (2019). Visualization in bayesian
 1055 workflow. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, *182*(2), 389–402.
- 1056 Gamer, M., Lemon, J., Gamer, M. M., Robinson, A., & Kendall's, W. (2012). Package “irr.” *Various*
 1057 *Coefficients of Interrater Reliability and Agreement*.
- 1058 Griffin, A. S., & Guez, D. (2014). Innovation and problem solving: A review of common mechanisms.
 1059 *Behavioural Processes*, *109*, 121–134. <https://doi.org/10.1016/j.beproc.2014.08.027>
- 1060 Griffin, A. S., Guez, D., Federspiel, I., Diquelou, M., & Lermite, F. (2016). Invading new environments:
 1061 A mechanistic framework linking motor diversity and cognition to establishment success. *Biological*
 1062 *Invasions and Animal Behaviour*, 26e46.
- 1063 Griffin, A. S., Guez, D., Lermite, F., & Patience, M. (2013). Tracking changing environments: Innovators
 1064 are fast, but not flexible learners. *PloS One*, *8*(12), e84907.

- 1065 Hadfield, J. (2010). MCMC methods for multi-response generalized linear mixed models: The MCMCglmm
 1066 r package. *Journal of Statistical Software*, 33(2), 1–22. <https://doi.org/10.18637/jss.v033.i02>
- 1067 Hadfield, J. (2014). *MCMCglmm course notes*. [http://cran.r-project.org/web/packages/MCMCglmm/
 1068 vignettes/CourseNotes.pdf](http://cran.r-project.org/web/packages/MCMCglmm/vignettes/CourseNotes.pdf)
- 1069 Hartig, F. (2019). *DHARMA: Residual diagnostics for hierarchical (multi-level / mixed) regression models*.
 1070 <http://florianhartig.github.io/DHARMA/>
- 1071 Hernán, M. A., & Robins, J. M. (2006). Instruments for causal inference: An epidemiologist’s dream?
 1072 *Epidemiology*, 360–372.
- 1073 Hlavac, M. (2018). *Stargazer: Well-formatted regression and summary statistics tables*. Central European
 1074 Labour Studies Institute (CELSI). <https://CRAN.R-project.org/package=stargazer>
- 1075 Homberg, J. R., Pattij, T., Janssen, M. C., Ronken, E., De Boer, S. F., Schoffelmeer, A. N., & Cuppen, E.
 1076 (2007). Serotonin transporter deficiency in rats improves inhibitory control but not behavioural flexibility.
 1077 *European Journal of Neuroscience*, 26(7), 2066–2073.
- 1078 Hutcheon, J. A., Chioloro, A., & Hanley, J. A. (2010). Random measurement error and regression dilution
 1079 bias. *Bmj*, 340, c2289. <https://doi.org/10.1136/bmj.c2289>
- 1080 Isden, J., Panayi, C., Dingle, C., & Madden, J. (2013). Performance in cognitive and problem-solving tasks
 1081 in male spotted bowerbirds does not correlate with mating success. *Animal Behaviour*, 86(4), 829–838.
- 1082 Kangas, B. D., & Bergman, J. (2017). Touchscreen technology in the study of cognition-related behavior.
 1083 *Behavioural Pharmacology*, 28(8), 623. <https://doi.org/10.1097/FBP.0000000000000356>
- 1084 Lajeunesse, M. J., Koricheva, J., Gurevitch, J., & Mengersen, K. (2013). Recovering missing or partial data
 1085 from studies: A survey of conversions and imputations for meta-analysis. *Handbook of Meta-Analysis in
 1086 Ecology and Evolution*, 195–206.
- 1087 Lea, S. E., Chow, P. K., Leaver, L. A., & McLaren, I. P. (2020). Behavioral flexibility: A review, a model,
 1088 and some exploratory tests. *Learning & Behavior*, 48(1), 173–187.
- 1089 Lefebvre, L., Whittle, P., Lascaris, E., & Finkelstein, A. (1997). Feeding innovations and forebrain size in
 1090 birds. *Animal Behaviour*, 53(3), 549–560. <https://doi.org/10.1006/anbe.1996.0330>
- 1091 Lin, G. (2020). *Reactable: Interactive data tables based on 'react table'*. [https://CRAN.R-project.org/
 1092 package=reactable](https://CRAN.R-project.org/package=reactable)
- 1093 Liu, Y., Day, L. B., Summers, K., & Burmeister, S. S. (2016). Learning to learn: Advanced behavioural
 1094 flexibility in a poison frog. *Animal Behaviour*, 111, 167–172.
- 1095 Logan, C. J. (2016). Behavioral flexibility in an invasive bird is independent of other behaviors. *PeerJ*, 4,
 1096 e2215.
- 1097 Logan, C. J., Avin, S., Boogert, N., Buskell, A., Cross, F. R., Currie, A., Jelbert, S., Lukas, D., Mares, R.,
 1098 Navarrete, A. F., et al. (2018). Beyond brain size: Uncovering the neural correlates of behavioral and
 1099 cognitive specialization. *Comparative Cognition & Behavior Reviews*.
- 1100 Logan, C. J., McCune, K., MacPherson, M., Johnson-Ulrich, Z., Rowney, C., Seitz, B., Blaisdell, A., Deffner,
 1101 D., & Wascher, C. (2021). *Are the more flexible great-tailed grackles also better at behavioral inhibition?*
 1102 <https://doi.org/10.31234/osf.io/vpc39>
- 1103 Logan, C. J., Shaw, R., Lukas, D., & McCune, K. B. (2022). *How to succeed in human modified environments*.
 1104 <http://corinalogan.com/ManyIndividuals/mil.html>
- 1105 Logan, CJ, McCune, KB, Chen, N, & Lukas, D. (2020). Implementing a rapid geographic range expansion
 1106 - the role of behavior and habitat changes. *In Principle Acceptance by PCI Ecology of the Version on 6
 1107 Oct 2020*. <http://corinalogan.com/Preregistrations/gxpobbehaviorhabitat.html>
- 1108 Logan, C., Lukas, D., Blaisdell, A., Johnson-Ulrich, Z., MacPherson, M., Seitz, B., Sevchik, A., & McCune,
 1109 K. (2022). Data: Is behavioral flexibility manipulatable and, if so, does it improve flexibility and problem
 1110 solving in a new context? *Knowledge Network for Biocomplexity, Data package*. [https://doi.org/10.5063/
 1111 F1RJ4GX6](https://doi.org/10.5063/F1RJ4GX6)
- 1112 Lukas, D., McCune, K., Blaisdell, A., Johnson-Ulrich, Z., MacPherson, M., Seitz, B., Sevchik, A., & Logan,
 1113 C. (2022). Behavioral flexibility is manipulatable and it improves flexibility and problem solving in a new
 1114 context: Post-hoc analyses of the components of behavioral flexibility. *EcoEvoRxiv*. [https://doi.org/10.
 1115 32942/osf.io/4ycps](https://doi.org/10.32942/osf.io/4ycps)
- 1116 Manrique, H. M., Völter, C. J., & Call, J. (2013). Repeated innovation in great apes. *Animal Behaviour*,
 1117 85(1), 195–202. <https://doi.org/10.1016/j.anbehav.2012.10.026>
- 1118 McCune, K., Blaisdell, A., Johnson-Ulrich, Z., Lukas, D., MacPherson, M., Seitz, B., Sevchik, A., & Logan,

1119 C. (2022). Repeatability of performance within and across contexts measuring behavioral flexibility.
1120 *EcoEvoRxiv*. <https://doi.org/10.32942/osf.io/kevqp>

1121 McCune, KB, MacPherson, M, Rowney, C, Bergeron, L, Folsom, M, & Logan, C. (2019). Is behavioral flexi-
1122 bility linked with exploration, but not boldness, persistence, or motor diversity? *In Principle Acceptance*
1123 *by PCI Ecology of the Version on 27 Mar 2019*. http://corinalogan.com/Preregistrations/g_exploration.html

1124

1125 McElreath, R. (2016). *Statistical rethinking: A bayesian course with examples in r and stan*. CRC Press.
1126 <https://doi.org/10.1201/9781315372495>

1127 McElreath, R. (2018). *Statistical rethinking: A bayesian course with examples in r and stan*. Chapman;
1128 Hall/CRC.

1129 McElreath, R. (2020). *Rethinking: Statistical rethinking book package*.

1130 McInerney, R. E. (2010). Multi-armed bandit bayesian decision making. *Univ. Oxford, Oxford, Tech. Rep*.

1131 Mikhalevich, I., Powell, R., & Logan, C. (2017). Is behavioural flexibility evidence of cognitive complexity?
1132 How evolution can inform comparative cognition. *Interface Focus*, 7(3), 20160121. <https://doi.org/10.1098/rsfs.2016.0121>

1133

1134 O’Hara, M., Huber, L., & Gajdon, G. K. (2015). The advantage of objects over images in discrimination
1135 and reversal learning by kea, nestor notabilis. *Animal Behaviour*, 101, 51–60.

1136 R Core Team. (2017). *R: A language and environment for statistical computing*. R Foundation for Statistical
1137 Computing. <https://www.R-project.org>

1138 Rayburn-Reeves, R. M., Stagner, J. P., Kirk, C. R., & Zentall, T. R. (2013). Reversal learning in rats
1139 (*rattus norvegicus*) and pigeons (*columba livia*): Qualitative differences in behavioral flexibility. *Journal*
1140 *of Comparative Psychology*, 127(2), 202.

1141 Revelle, W. (2014). *Psych: Procedures for psychological, psychometric, and personality research*. North-
1142 western University, Evanston, Illinois, 165, 1–10.

1143 Revelle, W. (2017). *Psych: Procedures for psychological, psychometric, and personality research*. North-
1144 western University. <https://CRAN.R-project.org/package=psych>

1145 Sawa, K., Leising, K. J., & Blaisdell, A. P. (2005). Sensory preconditioning in spatial learning using a touch
1146 screen task in pigeons. *Journal of Experimental Psychology: Animal Behavior Processes*, 31(3), 368.

1147 Schusterman, R. J. (1962). Transfer effects of successive discrimination-reversal training in chimpanzees.
1148 *Science*, 137(3528), 422–423.

1149 Seitz, B. M., McCune, K., MacPherson, M., Bergeron, L., Blaisdell, A. P., & Logan, C. J. (2021). Using
1150 touchscreen equipped operant chambers to study animal cognition. Benefits, limitations, and advice.
1151 *PloS One*, 16(2), e0246446.

1152 Shaw, R. C., Boogert, N. J., Clayton, N. S., & Burns, K. C. (2015). Wild psychometrics: Evidence for
1153 ‘general’ cognitive performance in wild new zealand robins, *petroica longipes*. *Animal Behaviour*, 109,
1154 101–111.

1155 Sol, D., Duncan, R. P., Blackburn, T. M., Cassey, P., & Lefebvre, L. (2005). Big brains, enhanced cognition,
1156 and response of birds to novel environments. *Proceedings of the National Academy of Sciences of the*
1157 *United States of America*, 102(15), 5460–5465. <https://doi.org/10.1073/pnas.0408145102>

1158 Sol, D., & Lefebvre, L. (2000). Behavioural flexibility predicts invasion success in birds introduced to new
1159 zealand. *Oikos*, 90(3), 599–605. <https://doi.org/10.1034/j.1600-0706.2000.900317.x>

1160 Sol, D., Székely, T., Liker, A., & Lefebvre, L. (2007). Big-brained birds survive better in nature. *Proceedings*
1161 *of the Royal Society of London B: Biological Sciences*, 274(1611), 763–769.

1162 Sol, D., Timmermans, S., & Lefebvre, L. (2002). Behavioural flexibility and invasion success in birds. *Animal*
1163 *Behaviour*, 63(3), 495–502.

1164 Spence, K. W. (1936). The nature of discrimination learning in animals. *Psychological Review*, 43(5), 427.

1165 Stan Development Team. (2020). *RStan: The R interface to Stan*. <http://mc-stan.org/>

1166 Summers, J., Lukas, D., Logan, C., & Chen, N. (2022). The role of climate change and niche shifts in
1167 divergent range dynamics of a sister-species pair. *EcoEvoRxiv*. <https://doi.org/10.32942/osf.io/879pe>

1168 Ushey, K., Allaire, J., Wickham, H., & Ritchie, G. (2020). *Rstudioapi: Safely access the RStudio API*.
1169 <https://CRAN.R-project.org/package=rstudioapi>

1170 Warren, J. (1965). Primate learning in comparative perspective. *Behavior of Nonhuman Primates*, 1,
1171 249–281.

1172 Warren, J. (1966). Reversal learning and the formation of learning sets by cats and rhesus monkeys. *Journal*

1173 *of Comparative and Physiological Psychology*, 61(3), 421.

1174 Warren, J. M. (1965). The comparative psychology of learning. *Annual Review of Psychology*, 16(1), 95–118.

1175 Wehtje, W. (2003). The range expansion of the great-tailed grackle (*quiscalus mexicanus gmelin*) in north
1176 america since 1880. *Journal of Biogeography*, 30(10), 1593–1607. [https://doi.org/10.1046/j.1365-2699.](https://doi.org/10.1046/j.1365-2699.2003.00970.x)
1177 [2003.00970.x](https://doi.org/10.1046/j.1365-2699.2003.00970.x)

1178 Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York. [https://](https://ggplot2.tidyverse.org)
1179 ggplot2.tidyverse.org

1180 Wickham, H., François, R., Henry, L., & Müller, K. (2021). *Dplyr: A grammar of data manipulation*.
1181 <https://CRAN.R-project.org/package=dplyr>

1182 Wilke, C. (n.d.). Cowplot: Streamlined plot theme and plot annotations for “ggplot2.” R package version
1183 0.9. 2; 2017. URL [Htts://CRAN. R-Project. Org/Package= Cowplot](https://CRAN.R-project.org/Package=Cowplot).

1184 Wolf, J. E., Urbano, C. M., Ruprecht, C. M., & Leising, K. J. (2014). Need to train your rat? There is an
1185 app for that: A touchscreen behavioral evaluation system. *Behavior Research Methods*, 46(1), 206–214.

1186 Wright, T. F., Eberhard, J. R., Hobson, E. A., Avery, M. L., & Russello, M. A. (2010). Behavioral flexibility
1187 and species invasions: The adaptive flexibility hypothesis. *Ethology Ecology & Evolution*, 22(4), 393–404.

1188 Xie, Y. (2013). Knitr: A general-purpose package for dynamic report generation in r. *R Package Version*,
1189 1(7).

1190 Xie, Y. (2017). *Dynamic documents with r and knitr*. Chapman; Hall/CRC.

1191 Xie, Y. (2018). Knitr: A comprehensive tool for reproducible research in r. In *Implementing reproducible*
1192 *research* (pp. 3–31). Chapman; Hall/CRC.

1193 Xie, Y. (2019). *formatR: Format r code automatically*. <https://CRAN.R-project.org/package=formatR>

1194 Zhu, H. (2021). *kableExtra: Construct complex table with 'kable' and pipe syntax*. [https://CRAN.R-](https://CRAN.R-project.org/package=kableExtra)
1195 [project.org/package=kableExtra](https://CRAN.R-project.org/package=kableExtra)