

Bayesian reinforcement learning models reveal how great-tailed grackles improve their behavioral flexibility in serial reversal learning experiments.

Lukas D^{1*} McCune KB² Blaisdell AP³ Johnson-Ulrich Z²
MacPherson M² Seitz B³ Sevchik A⁴ Logan CJ^{1*}

2024-06-06

Open...  access  code  peer review  data

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

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

Aur lie Coulon (2019) Can context changes improve behavioral flexibility? Towards a better understanding of species adaptability to environmental changes. *Peer Community in Ecology*, 100019. [10.24072/pci.ecology.100019](https://doi.org/10.24072/pci.ecology.100019). Reviewers: Maxime Dahirel and Andrea Griffin

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

Post-study manuscript (submitted to PCI Ecology for post-study peer review on 3 Jan 2022, revised and resubmitted Feb 2024): [preprint at EcoEvoRxiv](#), [rmd with code at github](#)

Abstract

Environments can change suddenly and unpredictably, so animals might benefit from being able to flexibly adapt their behavior through learning new associations. Reversal learning experiments, where individuals initially learn that a reward is associated with a specific cue before the reward is switched to a different cue, thus forcing individuals to reverse their learned associations, have long been used to investigate differences in behavioral flexibility among individuals and species. Here, we apply and expand newly developed Bayesian reinforcement learning models to gain additional insights into how individuals might dynamically adapt their behavioral flexibility if they experience repeated reversals in which cue is associated with a reward. Using data from simulations and great tailed grackles (*Quiscalus mexicanus*), we find that two parameters, the association updating rate, which reflects how much individuals weigh the most recent information relative to previously learned associations, and the sensitivity to learned associations, which reflects whether individuals no longer explore alternative options after having formed associations, are sufficient to explain the different strategies individuals display during the experiment. Individuals gain rewards more consistently if they have a higher association updating rate, because they learned that cues are reliable and they therefore can gain the reward consistently during one phase. The sensitivities to learned associations plays a role for the

grackles who experienced a series of reversals, where individuals with lower sensitivities are better able to explore the alternative option after a switch. The grackles who experienced the serial reversal adapted their behavioral flexibility through two different strategies. Some individuals showed more exploration such that they can quickly change to the alternative option after a switch even if they continue to occasionally choose the unrewarded option. Others stick to the previously learned associations such that they take longer to change after a switch, but, once they have reversed their associations consistently, choose the correct option. These strategies the grackles exhibited at the end of the reversal learning experiment also relate to their performance on multi-option puzzle boxes where there are different behaviors required to access rewards. Grackles with intermediate strategies solved fewer options to access the rewards than grackles with either of the extreme strategies, and they took longer to attempt a new option. Our approach offers new insights into how individuals react to uncertainty and changes in their environment, in particular showing that they can adapt their behavioral flexibility in response to their experiences.

Introduction

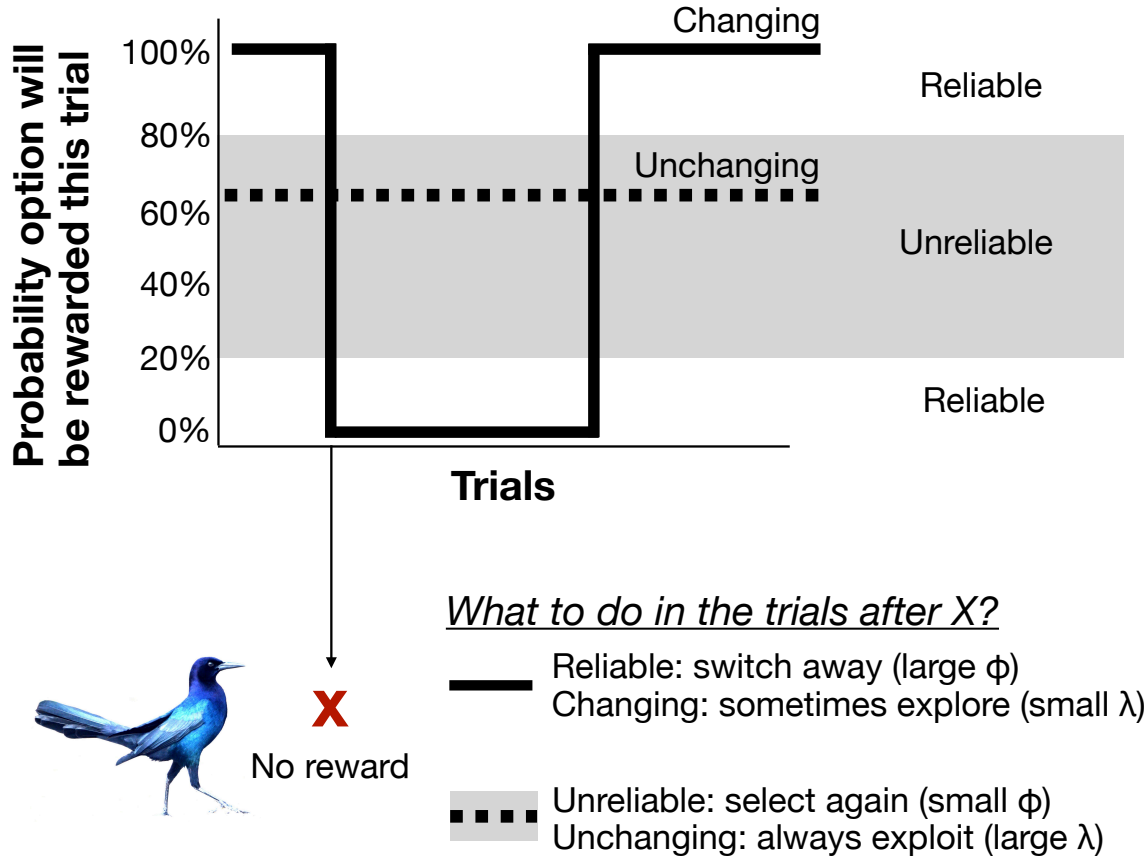
Most animals live in environments that undergo changes that can affect key components of their lives, such as where to find food or which areas are safe. Accordingly, individuals are expected to be able to react to these changes. One of the ways in which animals react to changes is through behavioral flexibility, the ability to change behavior when circumstances change (Shettleworth, 2010). The level of behavioral flexibility present in a given species is often assumed to have been shaped by selection, with past levels of change in the environment determining how well species might be able to cope with more rapidly changing (Sih, 2013) or novel environments (Sol et al., 2002). However, in another conception, behavioral flexibility is itself plastic (Wright et al., 2010). Behavioral flexibility arises because individuals update their information about the environment through personal experience and make that information available to other cognitive processes (Mikhalevich et al., 2017). Such modulation of behavioral flexibility is presumably relevant if the rate and extent of environmental change is variable and unpredictable Tello-Ramos et al. (2019). We are still limited in our understanding of when and how individuals might react to their experiences of environmental change.

Evidence that animals can change their behavioral flexibility based on their recent experience comes from serial reversal learning experiments. Serial reversal learning experiments have long been used to understand how individuals keep track of biologically important associations in changing environments Bitterman (1975). In these experiments, individuals are presented with multiple options associated with cues, such as different colors or locations, that differ in their reward. Individuals can repeatedly choose among the options to learn the associations between rewards and cues. After they show a clear preference for the most rewarded option, the rewards are reversed across cues, and individuals are observed to see how quickly they learn the changed associations. When they have reversed their preference, the reward is changed back to the other option, until the individual reverses their preference again, in a process called serial reversals. While the primary focus of these serial reversal learning experiments has been to measure differences in behavioral flexibility across individuals and species (Lea et al., 2020), several of these experiments show that behavioral flexibility is not a fixed trait, but that individuals can improve their performance if they experience repeated reversals Cauchoix et al. (2017). Here, we investigate how individuals might change their behavioral flexibility during serial reversal learning experiments to better understand what cognitive processes could lead to the observed differences and adjustments in behavioral flexibility Danwitz et al. (2022).

We recently found that great-tailed grackles (*Quiscalus mexicanus*; hereafter grackles) can be trained to improve how quickly they learn to change associations in a serial reversal learning experiment (Logan et al., 2023a). The experiment consisted of initially presenting birds with a light gray and a darkgray tube, only one of which contained a reward. After individuals chose one of the tubes, thus experiencing whether this color was either rewarded or not, the experiment was reset, with the reward being in the same colored tube as before. Once an individual chose the rewarded color more than expected by chance (passing criterion of choosing correctly in at least 17 out of the last 20 trials, based on a chi-square test), the reward was switched to the other color. Again, individuals made choices until they chose the now rewarded tube above the passing criterion. For one set of individuals, the trained group, we repeated the reversals of rewards from one color to the other until the birds reached the serial reversal passing criterion of forming a preference in 50 trials

or less in two consecutive reversals. The median number of trials birds in this trained group needed to reach the passing criterion during their first reversal was 75, which improved to 40 trials during their final reversal. Importantly, we found that, in comparison to a control group who only experienced a single reversal, trained grackles who experienced serial reversals also showed increased behavioral flexibility in other contexts. In particular, trained grackles performed better on multi-option puzzle boxes than control grackles, being faster to switch to a new access option on a box if the previous option was closed, and they solved more of the available access options (Logan et al., 2023a). This indicates that individuals did not just learn an abstract rule about the serial reversal learning experiment, but rather changed their overall behavioral flexibility in response to their experience.

A number of theoretical models have recently been developed that appear to reflect the potential cognitive processes individuals seem to rely on when making choices in reversal learning experiments (for a recent review see, for example, Frömer & Nassar, 2023). These theoretical models have been translated into Bayesian reinforcement learning models that can infer these processes from the data of individuals making choices during reversal learning Danwitz et al. (2022). The behavior of individuals in a reversal learning task is assumed to be guided by two cognitive processes Bartolo & Averbeck (2020). The first process reflects the *rate of updating associations* (which we refer to hereafter as ϕ , the Greek letter phi), the learning about the associations between the cues and potential rewards (or dangers). In the reinforcement learning models, this rate is reflected by the Rescorla-Wagner rule (Rescorla & Wagner, 1972). The rate weights the most recent information proportionally to the previously accumulated information for that cue (as a proportion, the rate can range between 0 and 1, see Equations 1-3). Individuals are expected to show different rates in different environments, particularly in response to the reliability of the cues (Figure 1). Lower updating rates are expected when associations are not perfect such that a single absence of a reward might be an error rather than indicating a new association. Higher updating rates are expected when associations are reliable such that individuals should update their associations quickly when they encounter new information Breen & Deffner (2023). The second process, the *sensitivity to their learned associations* (which we hereafter refer to as λ , the Greek letter lambda) reflects how individuals, when presented with a set of cues, might decide between these alternative options based on their learned associations of the cues. In the reinforcement learning model, the decision between different options is reflected by relative probabilities Danwitz et al. (2022), where the sensitivity to learned associations modifies the relative difference in learned rewards to generate the probabilities of choosing each option. A value of zero means individuals do not pay attention to their learned associations, but choose randomly, whereas increasingly larger values mean that individuals show strong biases in choice as soon as there are small differences in their learned associations (see Equations 1-3). Individuals with larger sensitivities will quickly prefer the option that previously gave them the highest reward (or the lowest danger), while individuals with lower sensitivities will continue to explore alternative options. Sensitivities are expected to reflect the rate of change in the environment (Figure 1), with larger sensitivities occurring when environments are static such that individuals start to exploit any differences they recognise as soon as possible. Lower sensitivities are expected when changes are frequent, such that individuals continue to explore alternative options when conditions change Breen & Deffner (2023).



124

125 **Figure 1** In serial reversal learning experiments, associations are reliable, such that if an option is associated
 126 with a reward, it is rewarded during every trial (white background). However, the associations between
 127 options and the rewards change across trials (solid line). In such environments, individuals are expected to
 128 gain the most rewards if they update their associations quickly (large ϕ) to switch away from an option if
 129 it is no longer being rewarded, and if they have small sensitivities to their learned associations to continue
 130 to explore all options to check if associations have changed again (small λ). In contrast, in unchanging but
 131 unreliable environments, the probability that an option is rewarded stays constant across trials (dotted lines),
 132 but is closer to 50% (gray background). In such environments, individuals are expected to gain the most
 133 rewards if they build their associations as average across many trials (small ϕ), and have high sensitivities to
 134 learned associations to exploit the option with the highest association (large λ). Grackle photo credit (CC
 135 BY 4.0): Dieter Lukas.

136 Here, we applied and modified the Bayesian reinforcement learning models to data from our grackle research
 137 on behavioral flexibility to assess if and how the cognitive processes might have changed as individuals
 138 experienced the serial reversal learning experiment. We previously found that the model can predict the
 139 performance of grackles in a static reversal learning task with a single reversal of a color preference (Blaisdell
 140 et al., 2021a). Grackles experiencing the serial reversal learning experiment are expected to infer that
 141 associations can frequently change but that, before and after a change, cues reliably indicate whether a
 142 reward is present or not. Based on the theoretical models, we predict that, in response to this experience,
 143 individuals increase their association-updating rate because cues are highly reliable, such that they can
 144 switch their associations as soon as there is a change Breen & Deffner (2023). In addition, we predict that
 145 individuals reduce their sensitivity to the learned associations, because the option that is rewarded switches
 146 frequently, requiring individuals to explore alternative options Leimar et al. (2024). Given that reversals
 147 in the associations are not very frequent, we also expect some variation in individuals in whether they
 148 switch to the newly rewarded option because they find the reward quickly through continued exploration
 149 (somewhat lower λ and higher ϕ) or because they quickly move away from the option that is no longer

rewarded (somewhat higher λ and lower ϕ). To assess these predictions, we addressed the following six research questions. With the first two research questions, we determined the feasibility and validity of our approach using simulations. With the other four questions, we analyzed the grackle data to determine how the association-updating rate and the sensitivity to learned associations reflect the variation and changes in behavioral flexibility in grackles.

1) Are the Bayesian reinforcement learning models sufficiently sensitive to detect changes that occur across the limited number of serial reversals that individuals participated in?

We use simulated data to answer this question because it allows us to determine how to apply the Bayesian reinforcement learning models to recover the likely changes in the cognitive processes during serial reversal learning. Previous applications of the models always combined the full sample of observations, so it is not clear whether these models are sufficiently sensitive to detect the changes we are interested in.

2) Do simulations confirm that a strategy of high association-updating (ϕ) and low sensitivity to learned associations (λ) is best to reduce errors in the serial reversal learning experiment?

We used simulations to systematically vary ϕ and λ to determine how the interaction of the two processes determine the behavior of individuals throughout the serial reversal learning experiment. We assessed whether, under the specific conditions in these experiments, information is reliable and changes occur frequently, and therefore the best strategy for individuals is to show high ϕ and low λ .

3) Which of the two parameters ϕ or λ explains more of the variation in the reversal learning experiment performance of the tested grackles?

Across both the trained (experienced serial reversals) and control (experienced single reversal) grackles, we assessed whether variation in the number of trials an individual needs to reach the criterion in a given reversal is better explained by their inferred association updating rate or by their sensitivity to learned associations.

4) Do the grackles who improved their performance through the serial reversal experiment show the predicted changes in ϕ and λ ?

If individuals learn the contingencies of the serial reversal experiment, they should be reducing their sensitivity to learned associations λ to explore the alternative option when rewards change, and increase their association-updating rate ϕ to quickly exploit the new reliably rewarded option.

5) Are some individuals better than others at adapting to the serial reversals?

In previous work, we found that there are individual differences that persist throughout the experiment, with individuals who required fewer trials to solve the initial reversal also requiring fewer trials in the final reversal after their training (McCune et al., 2023). We could expect that these individual differences are guided by consistency in how individuals solve the reversal learning paradigm, meaning they are reflected in individual consistency in ϕ and λ that persist through the serial reversals. In addition, it is not clear whether some grackles change their behavior more than others: for example, it could be that individuals who have a higher association-updating rate ϕ at the beginning of the experiment might also be better able to quickly change their behavior to match the particular conditions of the serial reversal learning experiment. Therefore, we also analyze whether the ϕ and λ values of individuals at the beginning predict how much they changed throughout the serial reversal learning experiment. Alternatively, given that the prediction for which sensitivity to learned association is best during a reversal (high sensitivity to stick to the learned associations) is different from the prediction for what is best right after a reversal (low sensitivity to explore the alternative option), the individuals who improved the most might end up with different strategies.

6) Can the ϕ or λ from the performance of the grackles during their final reversal predict variation in the performance on the multi-option puzzle boxes?

With the multi-option puzzle boxes, grackles would be expected to solve more options if they quickly update their previously learned associations when a previous option becomes unavailable (high ϕ). Given that, in the puzzle box experiment, individuals only receive a reward at any given option a few times, instead of repeatedly as during the reversal learning task, we predict that those individuals who are less sensitive to previously learned associations and instead continue to explore alternative options (low λ) can also gain more rewards.

Materials and Methods

Data

For questions 1 and 2, we re-analyzed data we previously simulated for power analyses to estimate sample sizes for population comparisons (Logan et al., 2023c). In brief, we simulated 20 individuals each from 32 different populations (640 individuals). The ϕ and λ values for each individual were drawn from a distribution representing that population, with different mean ϕ (8 different means) and mean λ (4 different values) for each population (32 populations as the combination of each ϕ and λ). The range for ϕ and λ values assigned to the artificial individuals in the simulations were based on the previous analysis of the single reversal data from grackles in a different population (Santa Barbara, California, USA, Blaisdell et al. (2021a)) to reflect the likely expected behavior. Based on their assigned ϕ and λ values, each individual was simulated to pass first through the initial association learning phase and, after they reached criterion (chose the correct option 17 out of the last 20 times), the rewarded option switched and simulated individuals went through the reversal learning phase until they again reached criterion. Each choice that each individual made was simulated consecutively, updating their internal associations with the two options based on their ϕ values and setting the probability of their next choice based on how their λ value weighted their associations to the two options. We excluded simulated individuals from the further analyses if they did not reach criterion either during the initial association or the reversal within 300 trials, the maximum that was also set for the experiments with the grackles. For each simulated individual, we had their assigned ϕ and λ values, as well as the series of choices they made during the initial association and the first reversal learning period.

For questions 3-5, we re-analyzed data of the performance of great-tailed grackles in serial reversal learning and multi-option puzzle box experiments (Logan et al., 2023a). The data collection was based on our preregistration that received in principle acceptance at PCI Ecology (Coulon, 2023). All of the analyses reported here were not part of the original preregistration. The data we use here were published as part of the earlier article (Logan et al., 2023b) and are available at the Knowledge Network for Biocomplexity’s data repository: https://knb.ecoinformatics.org/view/corina_logan.84.42. In brief, great-tailed grackles were caught in the wild in Tempe, Arizona, USA for individual identification (colored leg bands in unique combinations), and brought temporarily into aviaries for testing, before being released back to the wild. After training individuals to gain food from a yellow-colored tube, individuals then participated in the reversal learning tasks. A subset of individuals was part of the control group, where they learned the association of the reward with one color before experiencing one reversal to learn that the other color is rewarded (initial reward option was randomly assigned to either a dark-gray or a light-gray tube). The rewarded option was switched when grackles passed the criterion of choosing the rewarded option during 17 of the most recent 20 trials. This criterion was set based on earlier serial reversal learning studies, and is based on the chi-square test which indicates that 17 out of 20 represents a significant association. With this criterion, individuals can be assumed to have learned the association between the cue and the reward (Logan et al., 2022) rather than having randomly chosen one option more than the other. After their single reversal, the 11 control grackles participated in a number of trials with two identically colored tubes (yellow) which both contained a reward. This matched their general experiment participation to that of the trained group. The other subset of 8 individuals in the trained group went through a series of reversals until they reached the criterion of having formed an association (17 out of 20 choices correct) in less than 50 trials in two consecutive reversals. The individuals in the trained group needed between 6-8 reversals to consistently reach this threshold, with the number of reversals not being linked to their performance at the beginning or at the end of the experiment. After the individuals had completed the reversal learning experiment, both the control and trained individuals were provided access to two multi-access puzzle boxes, one made of wood and one made of plastic. The two boxes were designed with slight differences to explore how general the performance of the grackles was. The wooden box was made from a natural log, so was more representative of something the grackles might encounter in the wild. In addition, while both boxes had 4 possible ways (options) to access food, the four options on the wooden box were distinct compartments, each containing rewards, while the four options on the plastic box all led to the same reward. Grackles were tested sequentially on both boxes, where individuals could initially explore all options. After proficiency at an option was achieved (gaining food from this locus three times in a row), this option became non-functional by closing access to the option, and then the latency of the grackle to switch to attempting a different option was measured. If

they again successfully solved another option, this second option was also made non-functional, and so on. The outcome measures for each individual with each box were the average latency it took to switch to a new option and the total number of options they successfully solved.

The Bayesian reinforcement learning model

We used the version of the Bayesian model that was developed in Blaisdell et al. (2021b) and modified in Logan et al. (2023c) (see their Analysis Plan > “Flexibility analysis” for model specifications and validation). This model uses data from every trial of reversal learning (rather than only using the total number of trials to pass criterion) and represents behavioral flexibility using two parameters: the association-updating rate (ϕ) and the sensitivity to learned associations (λ). The model transforms the series of choices each grackle made based on two equations to estimate the most likely ϕ and λ that generated the observed behavior.

Equation 1 (attraction and ϕ): $A_{b,o,t+1} = (1 - \phi_b)A_{b,o,t} + \phi_b \pi_{b,o,t}$.

Equation 1 estimates how the associations A that individual b forms between the two different options (o_1 and o_2) and their expected rewards change from one trial to the next (time $t+1$) as a function of their previously formed associations $A_{b,o,t}$ (how preferable option o is to grackle b at time t) and recently experienced payoff π (in our case, $\pi = 1$ when they chose the correct option and received a reward in a given trial, and 0 when they chose the unrewarded option). The parameter ϕ_b modifies how much individual b updates its associations based on its most recent experience. The higher the value of ϕ_b , the faster the individual updates its associations, paying more attention to recent experiences, whereas when ϕ_b is lower, a grackle’s associations reflect averages across many trials. Association scores thus reflect the accumulated learning history up to this point. The association with the option that is not explored in a given trial remains unchanged. At the beginning of the experiment, we assume that individuals have the same low association between both options and rewards ($A_{b,1} = A_{b,2} = 0.1$).

Equation 2 (choice and λ): $P_{b,o,t+1} = \frac{\exp(\lambda_b A_{b,o,t})}{\sum_{o=1}^2 \exp(\lambda_b A_{b,o,t})}$.

Equation 2 expresses the probability P that an individual b chooses option o in the next trial, $t+1$, based on their learned associations of the two options with rewards. The parameter λ_b represents the sensitivity of a given grackle b to how different its associations to the two options are. As λ_b gets larger, choices become more deterministic and individuals consistently choose the option with the higher association even if associations are very similar. As λ_b gets smaller, choices become more exploratory, with individuals choosing randomly between the two options independently of their learned associations if λ_b is 0.

Equation 2 expresses the probability P that an individual b chooses option o in the next trial, $t+1$, based on the attractions. The parameter λ_b represents the rate of deviating from learned attractions of an individual. It controls how sensitive choices are to differences in attraction scores. As λ_b gets larger, choices become more deterministic and individuals consistently choose the option with the higher attraction even if attractions are very similar, as λ_b gets smaller, choices become more exploratory (random choice independent of the attractions if $\lambda_b=0$).

We implemented the Bayesian reinforcement learning model in the statistical language Stan (Stan Development Team, 2023), calling the model and analyzing its output in R (version 4.3.3) (R Core Team, 2023). The model takes the full series of choices individuals make (which of the two options did they choose, which option was rewarded, did they make the correct choice) across all their trials to find the ϕ and λ values that best fit these choices given the two equations: whether or not individuals chose the rewarded option was reflected as a categorical likelihood (yes or no) with probability P as estimated from equation 2, before updating the associations using equation 1. The model was fit across all choices, with individual ϕ and λ values estimated as varying effects. In the model, ϕ is estimated on the logit-scale to force the values to be positive before being converted back for equation 1 to update the associations, and λ is estimated on the log-scale to account for the exponentiation that occurs in equation 2. We set the priors for ϕ and λ to come from a normal distribution with a mean of zero and a standard deviation of one. We set the initial associations with both options for all individuals at the beginning of the experiment to 0.1 to indicate that they do not have an initial preference for either option but are likely to be somewhat curious about exploring

the tubes because they underwent habituation with a differently colored tube (see below). For estimations at the end of the serial reversal learning experiment, we set the association with the option that was rewarded before the switch to 0.7 and to the option that was previously not rewarded to 0.1. Note that when applying equation 1 in the context of the reversal learning experiment as most commonly used, where there are only rewards (positive association) or no rewards (zero association) but no punishment (negative association), associations can never reach zero because they change proportionally.

We used functions in the package “posterior” (Vehtari et al., 2021) to draw 4000 samples from the posterior (the default in the functions). We report the estimates for ϕ and λ for each individual (simulated or grackle) as the mean from these samples from the posterior. For the subsequent analyses where the estimated ϕ and λ values were response or predictor variables, we ran the analyses both with the single mean per individual as well as looping over the full 4000 samples from the posterior to reflect the uncertainty in the estimates. The analyses with the samples from the posterior provided the same estimates as the analyses with the single mean values, though with larger confidence estimates because of the increased uncertainty. In the results, we report the estimates from the analyses with the mean values. The estimates with the samples from the posterior can be found in the code in the rmd file at the repository. In analyses where ϕ and λ are predictor variables, we standardized the values that went into each analysis (either the means, or the respective samples from the posterior) by subtracting the average from each value and then dividing by the standard deviation. We did this to define the priors for the relationship on a more standard scale and to be able to more directly compare their respective influence on the outcome variable.

We also used the two equations analytically to more directly make predictions about how a specific ϕ and λ would influence the choices individuals make during the reversal learning. To derive the learning curves for individuals with different ϕ and λ , we incorporated the dynamic aspect of change over time by inserting the probabilities of choosing either the rewarded or the non-rewarded option from time $t-1$ as the likelihood for the changes in associations at time t .

Equation 3a (dynamic association for the rewarded option): $A_{r,t+1} = ((1-\phi) \times A_{r,t} + \phi \times \text{Reward}) \times P_t + (1-P_t) \times A_{r,t}$.

Equation 3b (dynamic association for the non-rewarded option): $A_{n,t+1} = (1-P_t) \times (1-\phi) \times A_{n,t} + P_t + (1-P_t) \times A_{n,t}$.

In equations 3a and 3b, the association with both the rewarded, A_r , and the non-rewarded, A_n , options change from time $t-1$ to time t depending on the association updating rate ϕ and the probability, P , that the association was chosen at time t . We used these equations to explore which combinations of ϕ and λ would lead to an individual choosing the rewarded option above the passing criterion within 50 trials after a switch in which option is rewarded. We assumed a serial reversal, and therefore set the initial associations after the presumed switch to 0.1 for the now rewarded option (previously unrewarded, so low association) and to 0.6 for the now unrewarded option (previously rewarded, so high association). For a given combination of ϕ and λ , we first used equation 2 to calculate the probability that an individual would choose the rewarded option during this first trial after the switch (where the remaining probability reflects the individual choosing the non-rewarded option), before using equations 3a and 3b to update the associations. We then repeated the calculations of the probabilities and the updates of the associations 50 times to determine whether individuals would reach the passing criterion with a given combination of ϕ and λ . For ϕ ranging between 0.02 and 0.10, we manually explored which λ would be needed.

1) Using simulations to determine whether the Bayesian serial reinforcement learning models have sufficient power to detect changes through the serial reversal learning experiment

We ran the Bayesian reinforcement learning model on these simulated data to understand the minimum number of choices per individual that would be necessary to recover the association-updating rate ϕ and the sensitivity to learned association λ values assigned to each individual.

To determine whether the Bayesian reinforcement learning model can accurately recover the simulated ϕ

and λ values from limited data, we applied the model first to only the choices from the initial association learning phase, next to only the choices from the first reversal learning phase, and finally from both phases combined. To estimate whether the Bayesian reinforcement learning model can recover the simulated ϕ and λ values without bias from either of the single or from the combined datasets, we correlated the estimated values with the values individuals were initially assigned:

$$\begin{aligned}\phi_{b,0} \text{ or } \lambda_{b,0} &\sim \text{Normal}(\mu, \sigma), \\ \mu &= \alpha + \beta \times \phi_{b,1} \text{ or } \lambda_{b,1}, \\ \alpha &\sim \text{Normal}(0, 0.1), \\ \beta &\sim \text{Normal}(1, 1), \\ \sigma &\sim \text{Exponential}(1),\end{aligned}$$

where a slope β between the assigned ($\phi_{b,0}$ or $\lambda_{b,0}$) and estimated ($\phi_{b,1}$ or $\lambda_{b,1}$) values close to 1 would indicate that the estimated values matched the assigned values.

This, and all following statistical models, were implemented using functions of the package ‘rethinking’ (McElreath, 2020) in R to estimate the association with stan. Following the social convention set in (McElreath, 2020), we report the mean estimate and the 89% compatibility interval from the posterior estimate from these models. For each model, we ran four chains with 10,000 iterations each (half of which were burn-in, and half samples for the posterior). We checked that the number of effective samples was sufficiently high and evenly distributed across parameters such that auto-correlation did not influence the estimates. We also confirmed that in all cases the Gelman-Rubin convergence diagnostic, \hat{R} , was 1.01 or smaller indicating that the chains had converged on the final estimates (Gelman & Rubin, 1995). In all cases, we also linked the model inferences back to the distribution of the raw data to confirm that the estimated predictions matched the observed patterns.

2) Using simulations to determine whether variation in ϕ or in λ has a stronger influence on the number of trials individuals might need to reach criterion in reversal learning experiments

We determined how the ϕ and λ values that were assigned to the simulated individuals influenced their performance in the reversal learning trials, building a regression model to determine which of the two parameters had a more direct influence on the number of trials individuals needed to reach criterion. We assumed that the number of trials followed a Poisson distribution because the number of trials to reach criterion is a count that is bounded at smaller numbers (individuals need at least 20 trials to reach the criterion), with a log-linear link, because we expect there are diminishing influences of further increases in ϕ or λ . The model is as follows:

$$\begin{aligned}v_{b,0} &\sim \text{Poisson}(\mu), \\ \log \mu &= \alpha + \beta_1 \times \phi + \beta_2 \times \lambda, \\ \alpha &\sim \text{Normal}(4.5, 1), \\ \beta_1 &\sim \text{Normal}(0, 1), \\ \beta_2 &\sim \text{Normal}(0, 1),\end{aligned}$$

where the prior for the intercept α was based on the average number of trials (90) grackles in Santa Barbara were observed to need to reach the criterion during the reversal (mean of 4.5 is equal to logarithm of 90, standard deviation set to 1 to constrain the estimate to the range observed across individuals). The priors for the relationships β_1 and β_2 with ϕ and λ were centered on zero, indicating that, *a priori*, we do not bias it toward a relationship.

3) Estimating ϕ and λ from the observed reversal learning performances of great-tailed grackles to determine which has more influence on variation in how many trials individuals needed to reach the passing criterion

We fit the Bayesian reinforcement learning model to the data of both the control and the trained grackles. Based on the simulation results indicating that the minimum sample required for accurate estimation are two learning phases across one switch (see below), we fit the model first to only the choices from the initial association learning phase and the first reversal learning phase for both control and trained individuals. For the control grackles, these estimated ϕ and λ values also reflect their behavioral flexibility at the end of the reversal learning experiment. For the trained grackles, we additionally calculated ϕ and λ separately for their final two reversals at the end of the serial reversals to infer the potential changes in the parameters. We fit the same regression model as with the simulated data to determine how ϕ and λ link to the number of trials grackles needed during their reversals.

4) Comparing ϕ and λ from the beginning and the end of the observed serial reversal learning performances to assess which changes more as grackles improve their performance

For the subset of grackles that were part of the serial reversal group, we calculated how much their ϕ and λ changed from their first to their last reversal. The model is as follows:

$$\begin{aligned} \phi \text{ or } \lambda &\sim \text{Normal}(\mu, \sigma), \\ \mu &= \alpha_b + \beta_b \times \text{reversal}, \end{aligned}$$

$$\begin{aligned} \begin{bmatrix} \alpha_b \\ \beta_b \end{bmatrix} &\sim \text{MVNormal}\left(\begin{bmatrix} \alpha \\ \beta \end{bmatrix}, S\right) \\ S &= \begin{pmatrix} \sigma_\alpha & 0 \\ 0 & \sigma_\beta \end{pmatrix} P \begin{pmatrix} \sigma_\alpha & 0 \\ 0 & \sigma_\beta \end{pmatrix} \end{aligned}$$

$$\begin{aligned} P &\sim \text{LKJcorr}(2), \\ \alpha &\sim \text{Normal}(5, 2), \\ \beta &\sim \text{Normal}(-1, 0.5), \\ \delta_b &\sim \text{Exponential}(1), \\ \sigma &\sim \text{Exponential}(1), \end{aligned}$$

where each grackle has two ϕ or λ values, one from the beginning (*reversal* equals 1) and one from the end of the serial reversal experiment (*reversal* equals 2). We assume that there are individual differences that persist through the experiment (intercept α_b) and that how much individuals change might also depend on their values at the beginning (multi-normal matrix correlation between the bird specific intercepts α_b and the bird specific changes between the reversals β_b).

We also fit a model to assess whether how much individuals improved in the number of trials from their first to their last reversal was linked more to their change in ϕ or to their change in λ . The model is as follows:

$$\begin{aligned} \Delta_b &\sim \text{Normal}(\mu, \sigma), \\ \mu &= \alpha + \beta_1 \times \Delta\phi_b + \beta_2 \times \Delta\lambda_b, \end{aligned}$$

$$\begin{aligned} \alpha_b &\sim \text{Normal}(40, 10), \\ \beta_1 &\sim \text{Normal}(0, 10), \\ \beta_2 &\sim \text{Normal}(0, 10), \\ \sigma &\sim \text{Exponential}(1), \end{aligned}$$

where Δ_b , the *improvement in the number of trials*, is the difference in the number of trials between the first and the last reversal and $\Delta\phi_b$ and $\Delta\lambda_b$ are the respective differences in these parameters between the beginning and the end of the serial reversal experiment.

5) Calculating whether individual differences in ϕ and λ persist throughout the serial reversal learning experiment and whether individuals differ in how much they change throughout the experiment

We checked whether the ϕ or λ values of individuals at the beginning (*first*) was associated with how much they changed (*change*, difference in values between beginning or end) or with the values they had at the end (*last*). The first part of the model is as follows:

$$\begin{aligned}\Delta\phi_b \text{ or } \Delta\lambda_b &\sim \text{Normal}(\mu, \sigma), \\ \mu_b &= \alpha + \beta \times \phi_{b,0} \text{ or } \lambda_{b,0}, \\ \alpha &\sim \text{Normal}(0,1), \\ \beta &\sim \text{Normal}(0,1), \\ \sigma &\sim \text{Exponential}(1),\end{aligned}$$

where $\phi_{b,0}$ and $\lambda_{b,0}$ are from the first reversal. The second part of the model is as follows:

$$\begin{aligned}\phi_{b,1} \text{ or } \lambda_{b,1} &\sim \text{Normal}(\mu, \sigma), \\ \mu &= \alpha + \beta \times \phi_{b,1} \text{ or } \lambda_{b,1}, \\ \alpha &\sim \text{Normal}(0,1), \\ \beta &\sim \text{Normal}(0,1), \\ \sigma &\sim \text{Exponential}(1),\end{aligned}$$

where $\phi_{b,1}$ and $\lambda_{b,1}$ are from the last reversal.

In addition, we assessed whether grackles at the end show the potential trade-off between ϕ and λ that could be expected in the serial reversal experiment. The model is as follows:

$$\begin{aligned}\phi_{b,1} &\sim \text{Normal}(\mu, \sigma), \\ \mu &= \alpha + \beta \times \lambda_{b,1}, \\ \alpha &\sim \text{Normal}(0,1), \\ \beta &\sim \text{Normal}(0,1), \\ \sigma &\sim \text{Exponential}(1),\end{aligned}$$

where

6) Linking ϕ and λ from the observed serial reversal learning performances to the performance on the multi-access boxes

We modified the models in the original article (Logan et al., 2023a) that linked performance on the serial reversal learning tasks to performance on the multi-access boxes, replacing the previously used independent variable of number of trials needed to reach criterion in the last reversal with the estimated ϕ and λ values from the last two reversals (trained grackles) or the initial discrimination and the first reversal (control grackles) (see below for explanation of these choices). With our expectation that ϕ and λ could be negatively correlated, we realized that grackles might be using different strategies when facing a situation in which cues change: some grackles might quickly discard previous information and rely on what they recently experienced (high ϕ and low λ), or they might rely on earlier information and continue to explore other options (low ϕ and high λ). Accordingly, we assumed that there also might be non-linear, U-shaped relationships between ϕ and/or λ and the performance on the multi-access box. For the number of options solved, we fit a binomial model with a logit link:

$$\begin{aligned}o_b &\sim \text{Binomial}(4, p), \\ \text{logit}(p) &\sim \alpha + \beta_1 \times \phi + \beta_2 \times \phi^2 + \beta_3 \times \lambda + \beta_4 \times \lambda^2,\end{aligned}$$

479 $\alpha \sim \text{Normal}(1, 1),$
 480 $\beta_1 \sim \text{Normal}(0, 1),$
 481 $\beta_2 \sim \text{Normal}(0, 1),$
 482 $\beta_3 \sim \text{Normal}(0, 1),$
 483 $\beta_4 \sim \text{Normal}(0, 1),$

484 where o_b is the number of options solved on the multi-access puzzle box, 4 is the total number of options, p is
 485 the probability of solving any one option across the whole experiment, α is the intercept, β_1 is the expected
 486 linear amount of change in o_b for every one unit change in ϕ in the reversal learning experiments, β_2 is the
 487 expected non-linear amount of change in o_b for every one unit change in ϕ^2 , β_3 the expected linear amount
 488 of change for changes in λ , and β_4 is the expected non-linear amount of change for changes in λ^2 .

489 For the average latency to attempt a new option on the multi-access puzzle box as it relates to trials to
 490 reverse (both are measures of flexibility), we fit a Gamma-Poisson model with a log-link:

491 $n_b \sim \text{Gamma-Poisson}(m_i, s),$
 492 $\log(m_i) \sim \alpha + \beta_1 \times \phi + \beta_2 \times \phi^2 + \beta_3 \times \lambda + \beta_4 \times \lambda^2,$
 493 $\alpha \sim \text{Normal}(1, 1),$
 494 $\beta_1 \sim \text{Normal}(0, 1),$
 495 $\beta_2 \sim \text{Normal}(0, 1),$
 496 $\beta_3 \sim \text{Normal}(0, 1),$
 497 $\beta_4 \sim \text{Normal}(0, 1), s \sim \text{Exponential}(1),$

498 where n_b is the average latency in seconds to attempt a new option on the multi-access box, m_i is the rate
 499 (probability of attempting an option in each second) per grackle (grackles with a higher rate have a smaller
 500 latency), s is the dispersion of the rates across grackles, α is the intercept, β_1 is the expected linear amount
 501 of change in latency for every one unit change in ϕ , β_2 is the expected non-linear amount of change in
 502 latency for every one unit change in ϕ^2 , β_3 the expected linear amount of change for changes in λ , and β_4
 503 the expected non-linear amount of change for changes in λ^2 .

Results

1) Power of the Bayesian reinforcement learning model to detect short-term changes in the association-updating rate ϕ and the sensitivity to learned associations λ

Applying the Bayesian reinforcement learning model to simulated data from only a single phase (initial association or first reversal) revealed that, while the model recovered the differences among individuals, the estimated ϕ and λ values did not match those the individuals had been assigned (Figure 2 shows the relationship between the assigned and estimated ϕ values when estimated from only the first reversal as an illustration). We realized that ϕ and λ values were consistently shifted, with the Bayesian estimation adjusting both parameters towards the mean and away from extreme values. Simulated individuals who were assigned large λ values were estimated to have a smaller λ values but in turn estimated to have ϕ values such that they would reach criterion in a similar number of trials because while the model assumed that they were more exploratory the model also assumed that they updated their associations more quickly. Similarly, individuals with large assigned ϕ values were estimated to have smaller ϕ values, but in turn were estimated to have larger λ values than those λ they were assigned. The estimated ϕ values did reflect the assigned ϕ values (here and hereafter, we report the posterior mean of the association with the 89% compatibility interval: +0.15, +0.06 to +0.23, n=626 simulated individuals), and the estimated λ values reflected the assigned λ values (+0.58, +0.48 to +0.68, n=626 simulated individuals). However, because the estimation from a single reversal did not accurately recover large values for either parameter, large ϕ and λ values were underestimates of the assigned values. In addition, this shift means that, even though simulated individuals were assigned ϕ and λ values randomly from across all possible combinations, the estimated values showed a strong positive correlation as the model had to make up the shifts in estimates of one parameter through shifting the estimate of the other parameter (slope of the correlation between the estimated λ and estimated ϕ values: +505, +435 to +570, n=626 simulated individuals).

In contrast, when we combined data from across the initial discrimination learning and the first reversal, the model accurately recovered the ϕ and λ values that the simulated individuals had been assigned (ϕ : +0.96, +0.70 to +1.21, n=626 simulated individuals; λ : +0.98, +0.92 to +1.05, n=626 simulated individuals) (Figure 2). While different combinations of ϕ and λ could potentially explain the series of choices during a single phase (initial discrimination and single reversal), these different combinations lead to different assumptions about how an individual would behave right after a reversal when the reward is switched to the alternative option, making it possible to infer the assigned value when combining behavioral choices from two phases (initial learning plus first reversal, or two subsequent reversals).

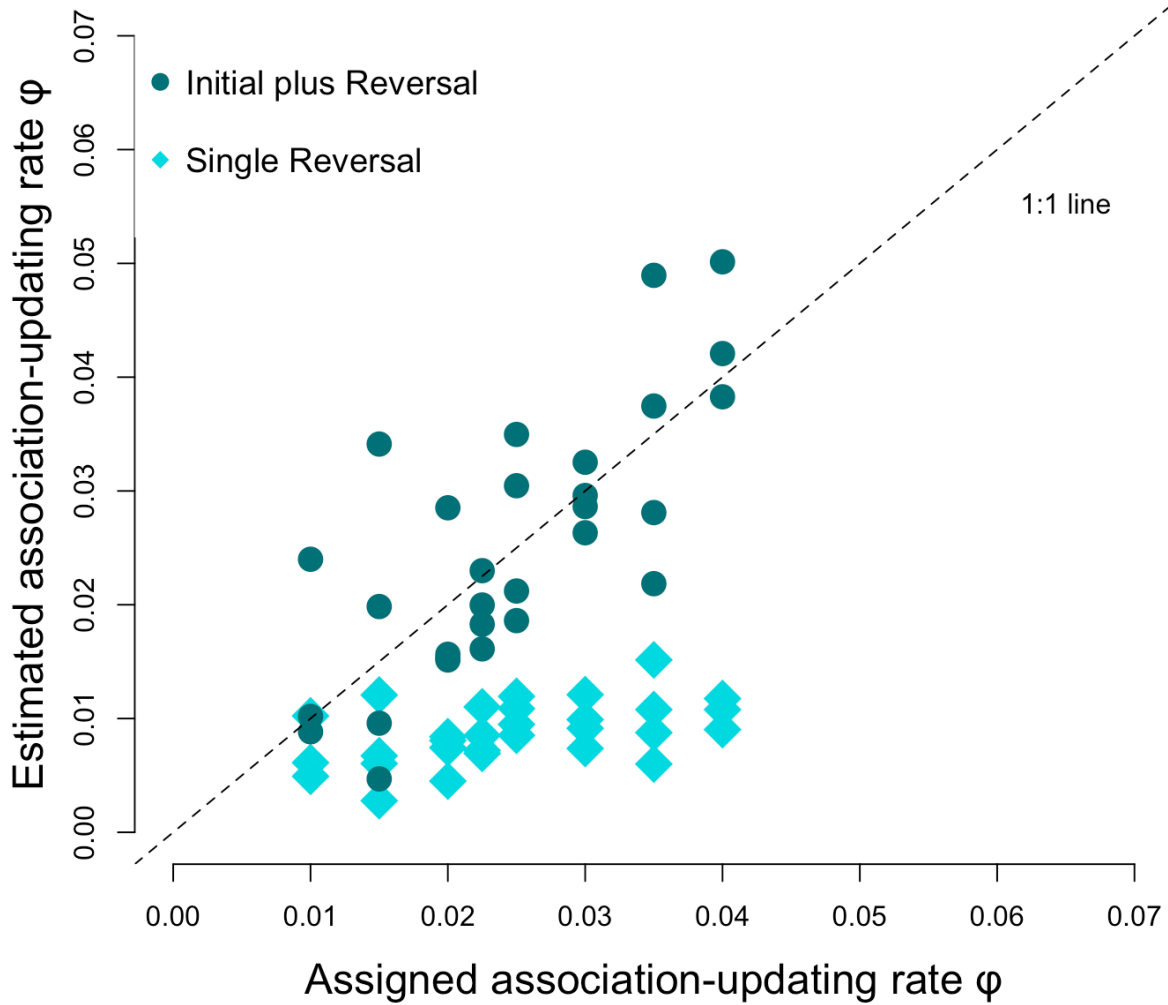


Figure 2: The ϕ values estimated by the model based on the choices made by 30 of the simulated individuals (y-axis) versus the ϕ values assigned to them (x-axis). Individuals were assigned a range of ϕ values, their choices were simulated and these values were used to back-estimate the ϕ . When ϕ was estimated based on the choices made only during the first reversal, the estimates were consistently lower than the assigned values, particularly for large ϕ values (lightblue squares). However, when ϕ was estimated based on the choices made during the initial association and the first reversal, the estimates were close to the assigned values (darkgreen circles). Patterns are similar for the relationship between the estimated and assigned λ values, and when ϕ and λ are estimated only from the trials during the initial association learning. Lines around the points indicate the compatibility intervals of the estimated values.

2) Predicted role of ϕ and λ on performance in the serial reversal learning task based on simulations

The ϕ values assigned to simulated individuals had a stronger influence on the number of trials they needed to pass the criterion during a reversal (-0.23, confidence interval: -0.24 to -0.23; n=626 simulated individuals) than their assigned λ values (-0.17, -0.18 to -0.16, n = 626 simulated individuals). In line with the prediction, there was a linear negative relationship between ϕ and the number of trials to reverse, with simulated

individuals needing fewer trials the more they updated their association based on their most recent experience. There also was, as predicted, an overall negative relationship between λ and the number of trials to reverse. Individuals generally needed few trials to reach the criterion if they were assigned a high λ value because they acted even on small differences in their learned associations. However, while individuals with small λ values can show large numbers of 150 or more trials to reach criterion because they are not sensitive to the differences in their learned associations, individuals with small λ values can also reach the criterion in small numbers of trials if they simultaneously quickly update their association because of their high ϕ values (Figure 3).

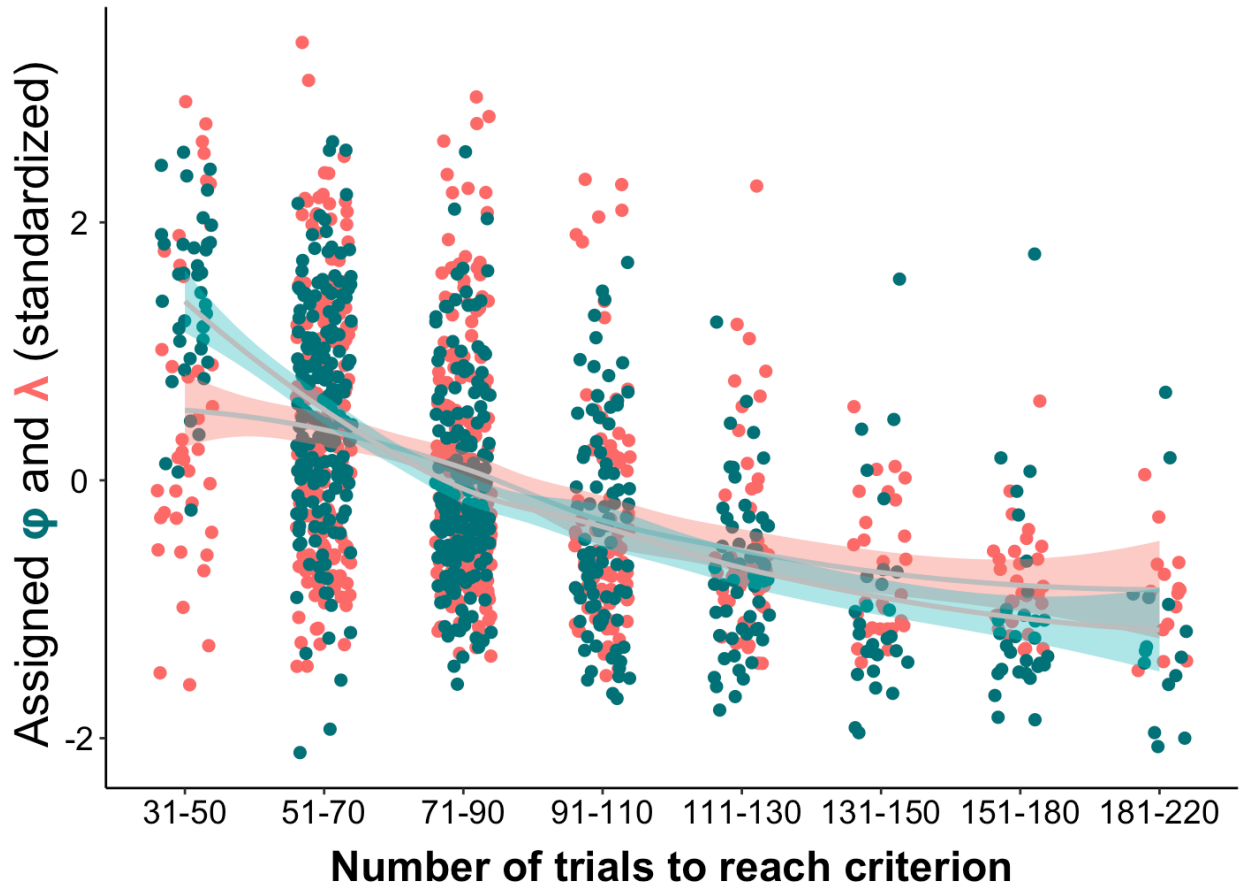
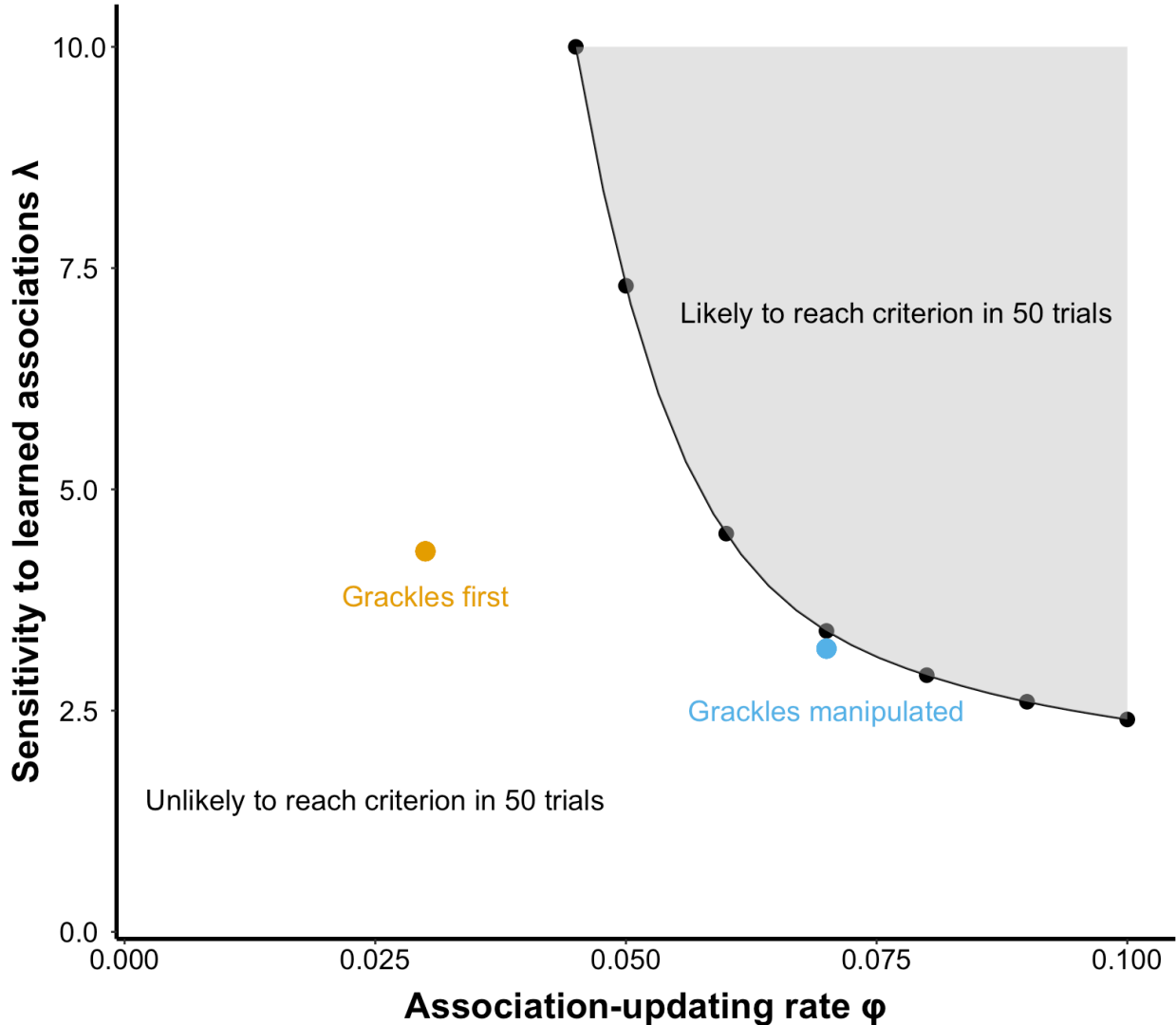


Figure 3. In the simulations, the ϕ values assigned to individuals (green) had a larger influence on the number of trials these individuals needed to reverse than their λ values (red). In general, individuals needed fewer trials to reverse if they had larger ϕ and λ values. However, relatively small λ values could be found across the range of reversal performances, whereas there was a more clear distinction with ϕ values (shaded lines represent compatibility intervals of the estimated relationship for these data). The number of trials to reach criterion are grouped into discrete blocks for easier illustration, but the analyses were performed on the raw values for each individual.

We performed an analytical assessment of this likely trade-off between the association updating rate ϕ and the sensitivity to the learned associations λ to identify the range of values we could expect in the serial reversal learning experiment. We assigned an hypothetical individual one of nine potential ϕ values in the range of 0.02 to 0.10 (steps differ by 0.01), assumed that this individual initially had the same association of the reward with both of the options (associations of 0.10 for light gray and 0.10 for dark gray), and assumed that this individual would choose each options 10 times during its first 20 trials. We calculated the associations to both options after the first 20 trials given the respective ϕ (e.g. with a ϕ of 0.10, the association with the rewarded option increases to 0.69 while the association with the unrewarded option declines to 0.03). Based on the differences in the two associations, we estimated the λ value necessary for

577 individuals to choose the rewarded option 85% in the next 20 trials (to reach the criterion of choosing the
 578 rewarded option in 17 out of 20 trials). We detected a clear negative, and exponential, trade-off between
 579 the necessary ϕ and λ values to reach the criterion (Figure 4): individuals with the highest ϕ value of 0.10
 580 only need a λ of 2.7 to reach the criterion, whereas individuals with a ϕ value of 0.02 need a λ of 9.5. This
 581 trade-off, where individuals can reach criterion during a reversal in few trials by either quickly updating their
 582 associations or by being highly sensitive to even small differences in their learned associations, means that in
 583 the serial reversal learning experiment individuals are expected to choose a strategy from across this range,
 584 and that doing so means they can also react to the sudden reversals in the reward location. In the serial
 585 reversal learning experiments, individuals will be able to reach the criterion more quickly during subsequent
 586 trials if they have, as predicted, a high ϕ and a low λ value. First, even if individuals were to choose randomly
 587 during the first trials after a reversal, individuals with a low ϕ need exponentially more trials to reverse their
 588 bias in associations between the two options. If an individual after one reversal has an association to the
 589 no longer rewarded option of 0.70 and to the now rewarded option of 0.10, with a ϕ of 0.02 it will take 48
 590 random trials until their association to the now rewarded options is higher than their association to the no
 591 longer rewarded option. In contrast, with a ϕ of 0.08 it will only take them 10 trials. Second, individuals
 592 with a high λ value will keep on choosing the previously rewarded option in almost all of their trials until
 593 this switch in associations occurs, further delaying the learning of the new associations. Individuals that
 594 have an association of 0.70 with the no longer rewarded option and 0.10 with the now rewarded option will
 595 choose the now rewarded option in 14% of cases if their λ is only 3, but only in 0.8% of cases if their λ is 8.



596

Figure 4. Individuals are more likely to reach the criterion of choosing the correct option 17 out of 20 times during the serial reversal trials if they update their associations quickly (high ϕ) and/or are sensitive to even small differences in their learned associations (high λ), because, during a reversal, recent information accurately predicts where the reward can be found. The figure shows this trade-off of individuals needing either high ϕ or high λ values to reach the criterion in a hypothetical situation where all individuals reach the criterion in 40 trials. This also means that if an individual has, for example, a high ϕ , their λ value becomes less important for reaching the criterion quickly. In this example, individuals with a ϕ of 0.10 will reach the criterion in 40 trials if their λ is at least 3.3. The figure also shows the median ϕ and λ values estimated for the grackles during their first reversal (yellow) when they needed about 70 trials to reach criterion and for the trained individuals during their last reversal (blue) when they did needed about 40 trials to reach criterion. During the training, grackles increased their ϕ to become efficient at gaining the reward and reaching the criterion, despite the concordant decline in λ .

3) Observed role of ϕ and λ on performance of grackles in the reversal learning task

For the grackles, we estimated ϕ and λ after the first reversal for all individuals, and additionally after the final reversal for the individuals who experienced the serial reversal learning experiment. The findings from the simulated data indicated that λ and ϕ can only be estimated accurately when calculated across at least one switch. In the simulation, we could combine the performance of individuals during the initial learning with the first reversal to estimate the parameters because the behavior during those two phases in the simulations was determined in the same way by the ϕ and λ values that individuals were assigned. We determined that we can also combine the first two phases for the grackles, because we found that the performance of the great-tailed grackles during the initial learning and the first reversal learning is correlated (+1.61, +1.53 to +1.69, n=19 grackles), with grackles needing about 28 trials more to reach criterion during the first reversal than they needed during the initial association learning. Therefore, we estimated ϕ and λ for the great-tailed grackles based on their performance in the initial discrimination plus first reversal, and for the trained grackles additionally based on their performance in the final two reversals. The inferred ϕ values for the grackles in Arizona range between 0.01 and 0.10, and the λ values between 2.1 and 6.5 (Figure 5).

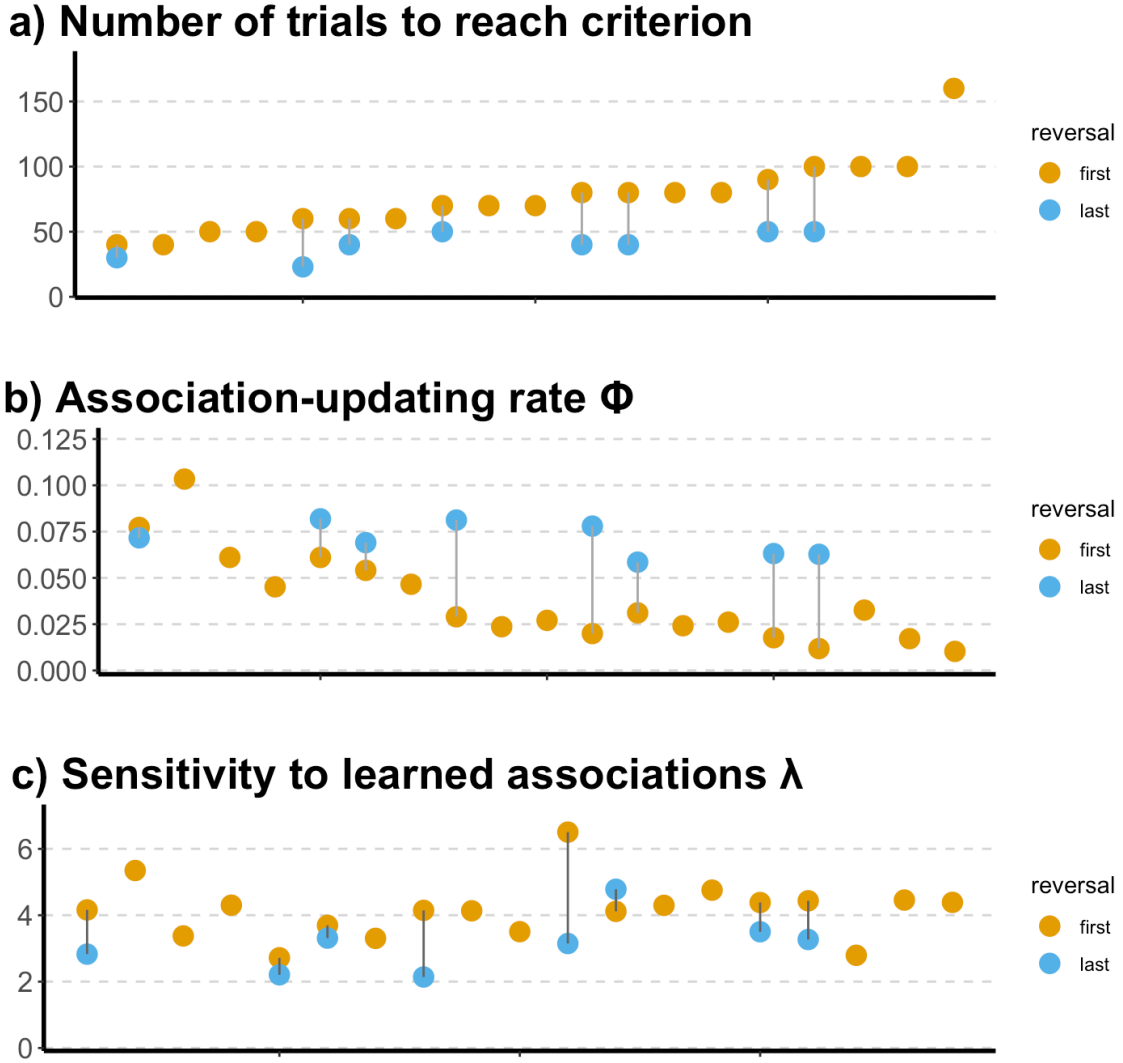


Figure 5. Comparisons of the different measures of ability in the reversal task for each of the 19 great-tailed grackles. The figure shows a) the number of trials to pass criterion for the first reversal (orange; all grackles) and the last reversal (blue; only trained grackles); b) the ϕ values reflecting the rate of updating associations with the two options inferred from the initial discrimination and first reversal (orange; all grackles) and from the last two reversals (blue; trained grackles); and c) the λ values reflecting the sensitivity to the learned associations inferred from the initial discrimination and first reversal (orange; all grackles) and from the last two reversals (blue; trained grackles). Individual grackles have the same position along the x-axis in all three panels. Grackles that needed fewer trials to reverse their preference generally had higher ϕ values, whereas λ appeared unrelated to the number of trials grackles needed during the first reversal. For the trained grackles, their ϕ values changed more consistently than their λ values, and the ϕ values of the trained individuals were generally higher than those observed in the control individuals, while their λ values remained within the range observed in the control group.

For the 19 grackles that finished the initial learning and the first reversal, only their ϕ (-20.69, -26.17 to -15.13; $n=19$ grackles), but not their λ (-0.22, -5.66 to 5.26, $n=19$ grackles), predicted the number of trials they needed to reach criterion during their first reversal (Figure 6). A grackle with a 0.01 higher ϕ than another individual needed about 10 fewer trials to reach the criterion. The slope between ϕ and the number of trials for the grackles was essentially identical to that observed in the simulations (-21.21 vs -20.48, Figure 6). The number of trials grackles needed to reach the criterion given their ϕ values fell right into the range

observed in the relationship between the ϕ and the number of trials observed among the simulated individuals (Figure 6). Even though the 8 trained grackles also appeared to need slightly fewer trials to reach criterion in their final two reversals if they had a higher ϕ , the limited variation in the number of trials and in ϕ and λ values among individuals means that there is no clear association (ϕ : -7.38, -15.97 to 1.28; λ : -4.00, 12.53 to 4.61, n=8 grackles).

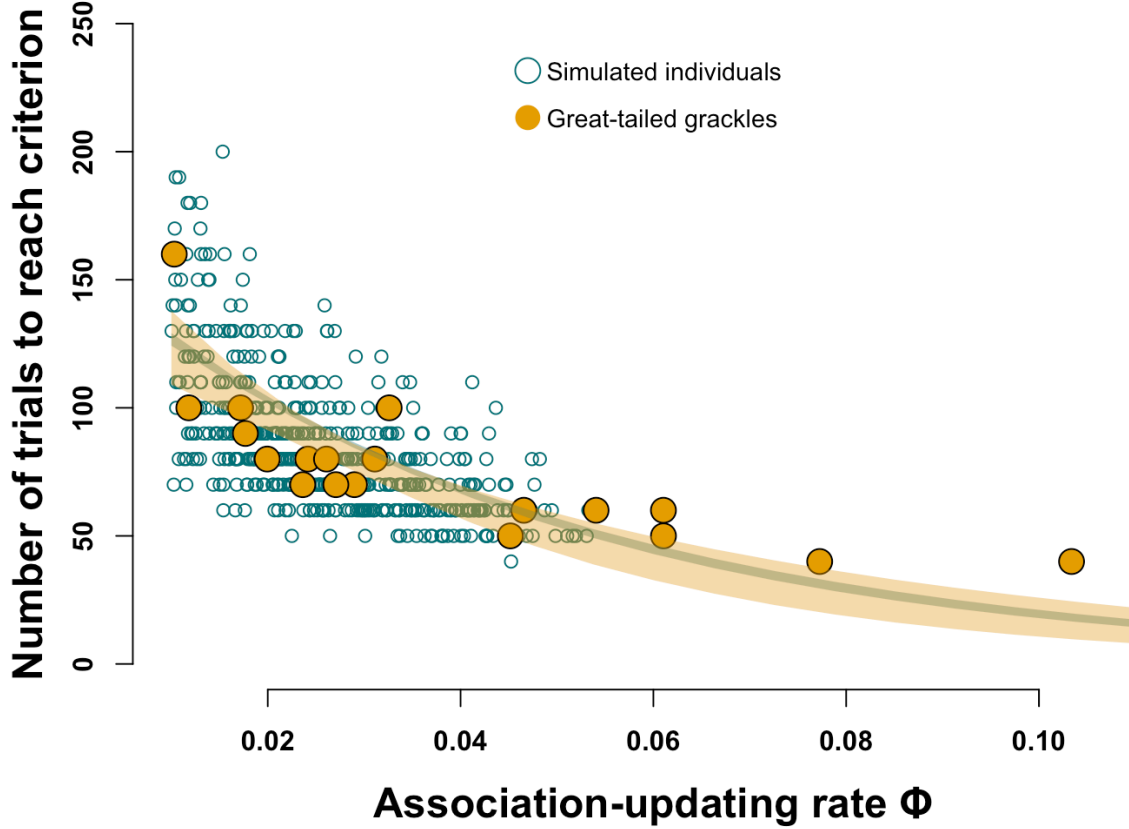


Figure 6. Relationship between ϕ and the number of trials grackles (yellow points) and simulated individuals (green circles) needed to reach criterion in their first trial. The observed grackle data falls within the range of the number of trials individuals with a given ϕ value are expected to need, and shows the same negative correlation between their ϕ and the number of trials as the simulated individuals (lines display the compatibility interval of the estimated relationships).

4) Changes in ϕ and λ through the serial reversal learning task

Great-tailed grackles who experienced the serial reversal learning reduced the number of trials they needed to reach the criterion from an average of 75 to an average of 40 (-30.02, -36.05 to -24.16, n=8 grackles). For the trained grackles, the estimated ϕ values more than doubled from 0.03 in their initial discrimination and first reversal (which is identical to the average observed among the control grackles who did not experience the serial reversals) to 0.07 in their last two reversals (+0.03, +0.02 to +0.05, n=8). The λ values of the trained grackles went slightly down from 4.2 (again, identical to control grackles) to 3.2 (-1.07, -1.63 to -0.56, n=8 grackles) (Figure 5). The values we observed after the training in the last reversal for the number of trials to reverse, as well as the ϕ and λ values estimated from the last reversal, all fall within the range of

variation we observed among the control grackles in their first and only reversal (Figure 5). This means that the training did not push grackles to new levels, but changed them within the boundaries of their natural abilities observed in the population.

As predicted, the increase in ϕ during the training fits with the outcome from the simulations: larger ϕ values were associated with fewer trials to reverse. The improvement the grackles showed in the number of trials they needed to reach the criterion from the first to the last reversal matched the changes in their ϕ values (+7.59, +1.54 to +14.22, n=8 grackles). The improvement did not match the change in their λ values (+2.17, -4.66 to 9.46, n=8 grackles), because, as predicted, the grackles in the training showed a decreased λ in their last reversal. This decrease in λ meant that grackles quickly found the rewarded option after a switch in which option was rewarded. In their first reversal grackles chose the newly rewarded option in 25% of the first 20 trials, in their final reversal the trained grackles chose correctly in 35% of the first 20 trials. Despite their low λ values, trained grackles still chose the rewarded option consistently because the increase in ϕ compensated for this reduced sensitivity (Figure 4; also see below).

5) Individual consistency in the serial reversal learning task

While we had previously found that differences among grackles in whether they needed many or few trials persisted through the serial reversals, we did not find similar consistency in either ϕ or λ . We found a negative correlation between the ϕ estimated from an individual's performance in the first reversal and how much their ϕ changed toward the value for their performance in the last reversal (-0.84, -1.14 to -0.52, n=8 grackles) such that individuals ended up with similar ϕ values to each other at the end of the training and their beginning and end ϕ values were not correlated (-0.21, -1.55 to 1.35, n=8 grackles). Similarly, individuals who started with a high λ changed more than individuals who already had lower a λ during the first reversal (-0.44, -0.76 to -0.10, n=8 grackles) and these changes were not consistent such that individual differences in λ did not remain through the serial reversal learning task (+0.17, -0.67 to +0.97, n=8 grackles). Individuals appeared to use different adjustments to their strategies to improve their performance through the training. There was a negative correlation between an individual's ϕ and λ after their last reversal (-0.39, -0.72 to -0.06, n=8 grackles), indicating that they ended up with different strategies from along the range of potential solutions. Some individuals quickly learn the new reward structure after a switch, but continue to explore the alternative option even after they have learned the new associations (high association-updating rate and low sensitivity to learned associations). Other individuals take longer to learn that the reward has switched but once they have reversed their associations they rarely choose the unrewarded option (Figure 7). Together, this suggests that all individuals improved by the same extent through the training such that the differences in their performances persisted, but they ended up with different strategies for how to quickly reach the criterion after a reversal by either having a high association updating rate or a low sensitivity to their learned associations.

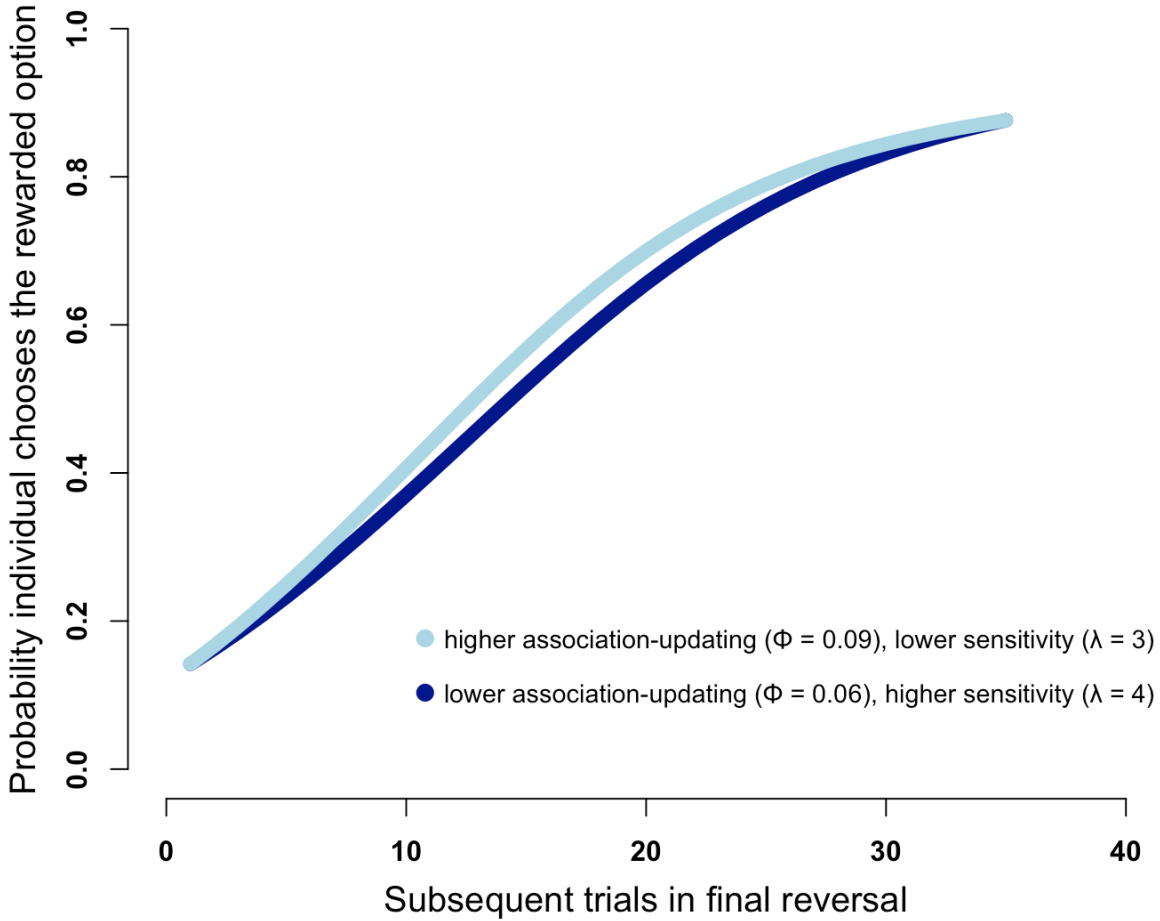


Figure 7. Predicted performance curves of individuals with different ϕ and λ values at the end of the serial reversal learning experiment based on the analytical formulas. We observed that, among the grackles who completed the serial reversal learning experiment, there was a negative correlation between their ϕ and λ , indicating that individuals used slightly different strategies to reach the criterion (choosing the rewarded option in 85% or more of trials) at equally few number of trials after the reward switched (when they had chosen the now rewarded option in 15% or less of trials). Individuals with a higher ϕ and lower λ (light blue line) quickly learn the new associations, but continue to explore the unrewarded option even after they have learned the association, leading to a curve with a more gradual increase throughout the trials. Individuals with a lower ϕ and higher λ (dark blue line) take longer to switch their associations, but once they do, they only rarely choose the non-rewarded option, leading to a more S-shaped curve where the initial increase in probability is lower and a more rapid rise later.

6) Association between ϕ and λ with performance on the multi-access boxes

We previously found that three measures of performance in the two multi-access puzzle boxes (number of options solved for both the wooden and the plastic multi-access puzzle box, latency to solve a new option on the plastic multi-access puzzle box) were correlated with the number of trials grackles needed to reach the criterion in the color tube reversal. We find that these measures also correlate with the underlying flexibility

parameters ϕ and λ . In particular, the number of options solved had a U-shaped association with the λ values individuals had at the end in their last reversal on both the plastic (estimate of association between number of options solved on plastic box and: ϕ : $= +0.03$, confidence interval -0.38 to $+0.43$; squared ϕ^2 : $= -0.16$, confidence interval -0.59 to $+0.28$; λ : $= +0.17$, confidence interval -0.27 to $+0.61$; squared λ^2 : $= +0.59$, confidence interval $+0.18$ to $+1.02$; $n=15$ grackles) and the wooden multi-access puzzle boxes (ϕ : -0.08 , -0.62 to $+0.47$; ϕ^2 : $+0.43$, -0.08 to $+0.97$; λ : $+0.03$, -0.50 to $+0.59$; λ^2 : $+0.63$, $+0.12$ to $+1.19$; $n=12$ grackles). Grackles who had either particularly low or particularly high sensitivities to their previously learned associations were more likely to solve all four options than grackles with intermediate values of λ (Figure 8). For the latency to attempt a new option on the plastic box there was also a U-shaped association, but with ϕ (ϕ : -0.66 , -1.30 to $+0.06$; ϕ^2 : $+0.58$, -0.06 to $+1.30$; λ : $+0.14$, -0.45 to $+0.70$; λ^2 : $+1.09$, $+0.28$ to $+1.87$; $n=11$ grackles) There was no association between the latency to attempt a new option on the wooden box with either ϕ (-0.62 , -1.46 to $+0.14$; ϕ^2 : $+0.39$, -0.47 to $+1.26$; $n=11$ grackles) nor λ ($+0.13$, -0.66 to $+0.86$; λ^2 : $+0.32$, -0.62 to $+1.35$; $n=11$ grackles). Grackles with either particularly high or particularly low rates of updating their associations took longer to attempt a new option than grackles with intermediate values of ϕ (Figure 8).

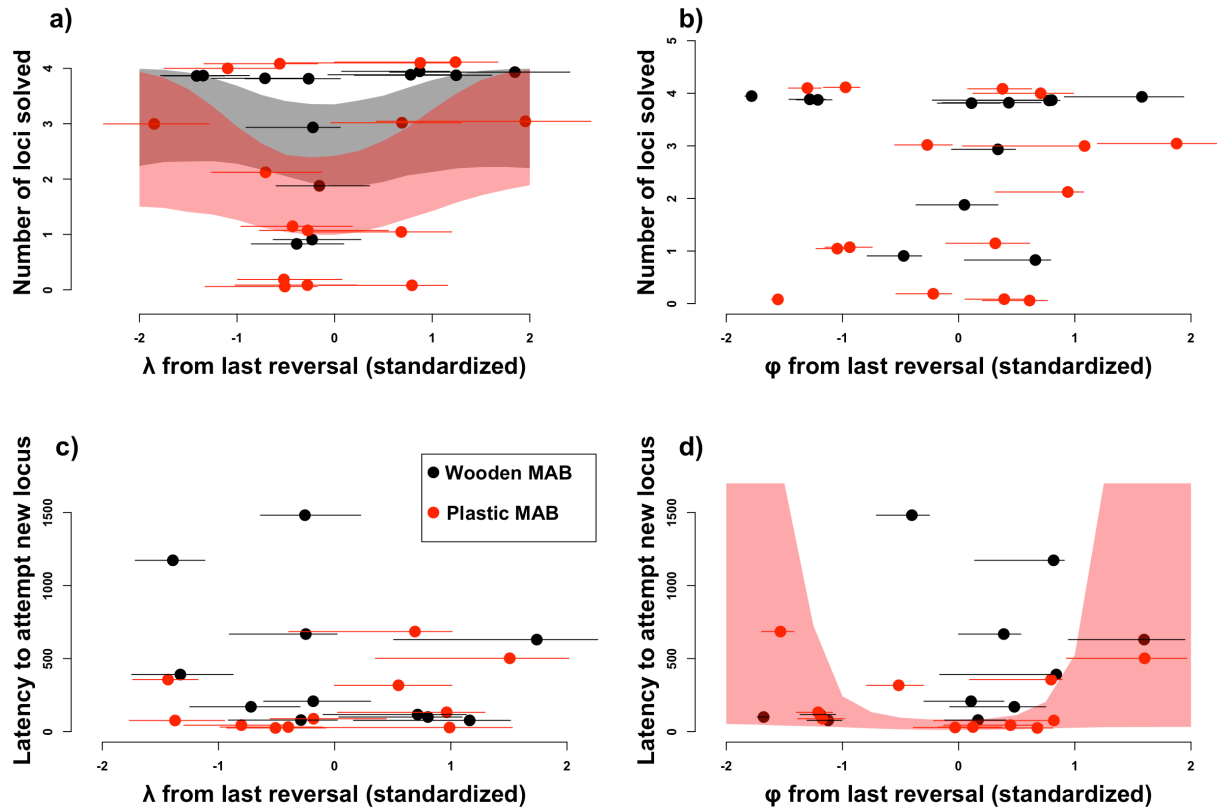


Figure 8. Relationships between ϕ and λ from the last reversal and performance on the wooden (black dots) and plastic (red dots) multi-access puzzle boxes. Grackles with intermediate λ values in their last reversal (a) were less likely to solve all four options on both boxes than grackles with either high or low λ values. Grackles with intermediate ϕ values have a shorter latency to attempt a new option on the plastic box (d). There are no clear relationships between ϕ and the number of options solved on either box (b), λ and the latency to attempt an option on either box (c), or (d) ϕ and the latency to attempt a new option on the wooden box. The ϕ and λ values change slightly between the top and bottom rows because the sample differs between boxes, and values were standardized for each plot.

Discussion

Our analyses indicate that applying a more mechanistic model to understand the behavior of great-tailed grackles in a serial reversal learning experiment can provide additional insights into the potential components of behavioral flexibility and their dynamic changes. First, the simulations showed that the Bayesian reinforcement learning model accurately captures variation in the behavior of individuals in the serial reversal learning experiment and that the two key parameters ϕ , the association-updating rate, and λ , the sensitivity to learned associations, can be reliably inferred if we combine at least two association learning periods across a switch in the rewarded options. This provides the opportunity to also infer whether and how individuals who experience the serial reversal learning experiment dynamically change their behavioral flexibility. Second, in line with our prediction, the simulations indicate that higher ϕ and lower λ mean that individuals should reach the reversal learning criterion in fewer trials. However, we observe that for a single reversal ϕ is more important and that λ simply sets a threshold on the number of trials individuals need to consistently choose the rewarded option. Third, post-hoc analyses of grackle serial reversal learning data revealed that, contrary to our prediction but in line with the simulation results, ϕ but not λ explained more of the interindividual variation in how many trials individuals needed to reach criterion during a reversal. Fourth, matching these observations, we found that the primary component of flexibility that was trained during the serial reversal experiments was ϕ , which more than doubled between the first and last reversals, whereas λ slightly declined, as expected based on the simulations. Fifth, while individual differences in performance persist across the serial reversals, the underlying changes in ϕ and λ are not predictable based on their initial values. Grackles appear to use different strategies to improve their performance during the serial reversal experiment, with some individuals showing more changes in their association-updating rate but less in their sensitivity to learned associations, while others show the opposite, leading to a negative correlation between the inferred ϕ and λ values among the individuals at the end of the serial reversal learning experiment. Finally, these different strategies to improve their behavioral flexibility that individuals revealed in the serial reversal learning experiment subsequently also influenced their behavior in a different experimental test of behavioral flexibility. Grackles with intermediate values of λ (and ϕ) solved fewer options on both multi-access puzzle boxes than grackles with either high or low λ (and low or high ϕ), and grackles with intermediate values of ϕ have shorter latencies to attempt a new option. Accordingly, the grackles appeared to react to the predictability of the associations and the frequent switches of the reward location that they experienced during the serial reversal learning experiment to adjust their behavioral flexibility.

Previous analyses of reversal learning performance in wild-caught animals have often focused on summaries of the choices individuals make (e.g. Bond et al., 2007), setting criteria to define success and how much individuals sample or explore the different options versus acquire or exploit the reward (e.g. Federspiel et al., 2017). These approaches are more descriptive, making it difficult to link the differences to specific processes and to predict how variation in behavior might transfer to other tasks. While there have been attempts to identify potential rules that individuals might learn during serial reversal learning (Spence, 1936; Warren, 1965a; Warren, 1965b; Minh Le et al., 2023), these rules were often about abstract switches to extreme strategies (e.g. win-stay / lose-shift) and therefore could not account for the full variation in the behavior. In contrast, the Bayesian reinforcement learning model with its two parameters of the association-updating rate and the sensitivity to learned associations has a clear theoretical foundation and appears to be sufficient to accurately represent the behavior of grackles in the serial reversal experiment. The previously described rules, including dramatic shifts in strategies, can be recovered with the dynamic Bayesian reinforcement learning model, including the different ‘learning curves’ that we observe among individuals (e.g. Gallistel et al., 2004). Applying the Bayesian reinforcement model to (serial) reversal learning experiments can provide several benefits to our understanding of behavioral flexibility. First, it highlights the key pieces of information that individuals likely pay attention to when adjusting their behavior. This provides ways to also link their performances and inferred cognitive abilities to how they experience and react to their natural environments. In particular, literature on foraging behavior that focuses on the likely trade-offs between the exploration versus exploitation of different options has a similar focus on gaining information (exploration) versus decision making (exploitation) (Kramer & Weary, 1991; Berger-Tal et al., 2014; Addicott et al., 2017). Having a mechanistic model for the behavioral choices can also help to design better and alternative experiments. Simulating the likely behavioral choices of individuals can help to decide how to track the progress of individuals and when to switch rewards (Logan et al., 2023a). Deciding on which external

conditions might matter most to a given group of individuals can help to determine which parameters to vary and can help to adapt the model further. For example, it has been extended to allow for unpredictability in the association between the cue and the reward (Gershman, 2018; Danwitz et al., 2022) or to assume that experiencing a reward will update the association more than not experiencing a reward (Metha et al., 2020). Our advance here was to make the model dynamic to determine how individuals adjust their behavior during the serial reversal learning experiment.

The dynamic model shows that behavioral flexibility in the grackles is not a fixed trait, but individuals can change their flexibility in response to their experiences. Grackles coming into the experiment already had different strategies, suggesting that they had different experiences of how predictable cues are and how frequently their environment changes. In general, the association-updating rate ϕ appears to explain more of the variation in how many trials individuals need to reach the criterion of consistently choosing the rewarded option during a single phase. The importance of the association-updating rate for the performance of the grackles in the reversal learning experiment matches what has been reported for squirrel monkeys (Bari et al., 2022). In contrast, the sensitivity to learned associations λ appears to set a threshold on the performance during a single phase, but appears more important as the rewards switch more frequently. In the serial reversal learning experiments, we observed an initial decline in performance, with most grackles needing more trials in the second and third reversal compared to the first, before improving and reaching the criterion in 50 trials or less (Logan et al., 2023a). This initial increase likely reflects that grackles need to distinguish between the absence of a reward at the previously rewarded location reflects stochastic variation in the association between the cue and the reward or an actual switch in reward structure. In a stochastic environment, individuals can gain more reward if they do not update their associations quickly, but stick with an option that previously gave them high rewards (Woo et al., 2023). In their natural environment, most cues are presumably not perfect such that their initial expectation might be that the particular tube just did not have a reward that time, but should still provide rewards frequently, thus explaining their initial decline in performance. Only after several switches is there sufficient information for the grackles to infer that the cues are highly reliable and the switches are relatively frequent. This is when they show the increase in their association-updating rate ϕ , which on average doubled across individuals, changing more for individuals who started off with lower ϕ values. Grackles also changed their sensitivity to the learned associations during the serial reversals, in line with the prediction that they benefit from being open to exploring the alternative option when the reward structure frequently switches.

Most animals that have been tested in serial reversal learning experiments thus far show improvements throughout the consecutive reversals, suggesting that most species can adapt their behavioral flexibility in response to the predictability and stability of their environments (e.g. Warren & Warren, 1962; Komischke et al., 2002; Bond et al., 2007; Strang & Sherry, 2014; Chow et al., 2015; Cauchoux et al., 2017; Degrande et al., 2022; Erdsack et al., 2022). For the grackles, the serial reversals pushed individuals to levels that were already observed in some individuals at the beginning of the experiment, meaning that the change within the experiment is within the natural range of abilities also observed in the wild. While there were individual differences in how individuals performed (McCune et al., 2023), all individuals changed depending on their experiences. Among the trained grackles, who all quickly switched to consistently gain the reward, we observed different strategies. On the one side, there are grackles who change gradually throughout an association phase, already choosing the newly rewarded option at the beginning but continuing to explore the alternative non-rewarded option throughout. These are the individuals with a high association-updating rate and low sensitivity to learned associations. On the other side are grackles who take longer to choose the newly rewarded option after a switch, but once they discover which option is rewarded, quickly reverse their preference. These are the individuals with low association-updating rates and high sensitivities to learned associations. With the variables we measured here, we could not predict which strategies grackles ended up with after the serial reversals. We observed additional strategies with different combinations of ϕ and λ across the grackles during their first reversal, but these are not efficient in the serial reversal learning experiment and instead are more suited to unpredictable and less frequently changing environments. How frequently and how quickly individuals change their behavioral flexibility in their natural environments is unclear. Individual differences might persist if their different behavioral flexibility leads them to continue to experience their environment differently. For the grackles, we have some indication that after releasing them back to their original environments, differences in behavioral flexibility between the trained and control

individuals persisted for at least several months, with individuals who had changed their ϕ and λ appearing to switch more frequently between food types and foraging techniques (Logan et al., 2024).

The analyses linking ϕ and λ to the performance on the multi-access boxes show that the different strategies grackles ended up with to improve their performance during the serial reversal learning experiment subsequently appeared to influence how they solved the multi-access box. The negative correlation between ϕ and λ prompted us to explore whether the relationship between these two variables and the performance on the multi-access boxes could be non-linear. We detected U-shaped relationships between ϕ and λ and how individuals performed on the multi-access puzzle boxes. First, grackles with intermediate ϕ values showed shorter latencies to attempt a new option. This could reflect that grackles with high ϕ values take longer because they formed very strong associations with the previously rewarded option, while grackles with small ϕ values take longer because they do not update their associations even though the first option is no longer rewarded or because they do not explore as much because of their small λ . Second, we found that grackles with intermediate values of λ solved fewer options. This could indicate that grackles with a small λ are more likely to explore new options while grackles with a large λ , and low ϕ are less likely to return to an option that is no longer rewarded. Given that there was also a positive correlation between the number of options solved and the latency to attempt a new options, there might be a trade-off, where grackles with extreme ϕ and λ values solve more options, but need more time, whereas grackles with intermediate values have shorter latencies, but solve fewer options. We are limited though in our interpretation by the small sample sizes. More detailed studies would be needed in order to fully understand how the association-updating rate and the sensitivity to learned associations might shape performance on the multi-access puzzle boxes. In addition, it is also possible that performance on the multi-access boxes relies on other cognitive abilities in which individuals may differ. For example, we previously found that grackles who are faster to complete an inhibition task, where they had to learn to not react to a cue in order to wait for a trial in which a different cue could result in gaining a reward, were slower to switch options on the boxes (Logan et al., 2021). As such, variation in self control may affect performance on flexibility and innovation tasks by decreasing exploratory behaviors. However, all these analyses are exploratory and based on a small sample, so these interpretations are speculative and further investigation is needed to understand how potential cognitive abilities shape performance on such tasks.

Overall, these findings indicate the potential benefits of applying more mechanistic models to psychological experiments. Inferring the cognitive processes potentially underlying behavior can allow us to make clearer predictions about how the performance in one experiment might translate to other paradigms and to behavior in the wild. For the serial reversal learning paradigm, we could expect that the previously observed differences in whether performance links with performance in other experiments like innovation or inhibition. For example, they correlate positively in gray squirrels (Chow et al., 2016), negatively in Indian mynas (Griffin et al., 2013), and positively, negatively, or not at all depending on the trait in great-tailed grackles (Logan, 2016 and this article). This variation could be linked to differences in whether the association-updating rate or the sensitivity to learned associations plays a larger role in the reversal performance in a given species and, in particular, for the other trait. The advanced capabilities of reflecting behavioral choices directly in a Bayesian framework offers an opportunity for the field of comparative cognition to implement more informed assessments of cognitive abilities and the factors shaping them.

Author contributions

Lukas: Hypothesis development, simulation development, data analyses, data interpretation, write up, revising/editing.

McCune: Added MAB log experiment, protocol development, data collection, revising/editing.

Blaisdell: Prediction revision, revising/editing.

Johnson-Ulrich: Data collection, revising/editing.

MacPherson: Data collection, revising/editing.

Seitz: Prediction revision, revising/editing.

Sevchik: Data collection, revising/editing.

Logan: Hypothesis development, protocol development, data collection, data analysis, data interpretation, revising/editing.

Funding

This research is funded by the Department of Human Behavior, Ecology and Culture at the Max Planck Institute for Evolutionary Anthropology.

Ethics

The research on the great-tailed grackles followed established ethical guidelines for the involvement and treatment of animals in experiments and received institutional approval prior to conducting the study (US Fish and Wildlife Service scientific collecting permit number MB76700A-0,1,2; US Geological Survey Bird Banding Laboratory federal bird banding permit number 23872; Arizona Game and Fish Department scientific collecting license number SP594338 [2017], SP606267 [2018], and SP639866 [2019]; California Department of Fish and Wildlife scientific collecting permit number S-192100001-19210-001; Institutional Animal Care and Use Committee at Arizona State University protocol number 17-1594R; Institutional Animal Care and Use Committee at the University of California Santa Barbara protocol number 958; University of Cambridge ethical review process non-regulated use of animals in scientific procedures: zoo4/17 [2017]).

Conflict of interest disclosure

We, the authors, declare that we have no financial conflicts of interest with the content of this article. CJ Logan is a Recommender and, until 2022, was on the Managing Board at PCI Ecology. D Lukas is a Recommender at PCI Ecology.

Acknowledgements

We thank our PCI Ecology recommender, Aurelie Coulon, and reviewers, Maxime Dahirel and Andrea Griffin, for their feedback on this preregistration; and the reviewers of this manuscript for their constructive feedback that helped with the framing of the study; Julia Cissewski for tirelessly solving problems involving financial transactions and contracts; Sophie Kaube for logistical support; and Richard McElreath for project support.

References

Addicott MA, Pearson JM, Sweitzer MM, Barack DL, Platt ML (2017) A primer on foraging and the explore/exploit trade-off for psychiatry research. *Neuropsychopharmacology*, **42**, 1931–1939. <https://doi.org/10.1038/npp.2017.108>

- Agrawal S, Goyal N (2012) Analysis of thompson sampling for the multi-armed bandit problem. In: *Conference on learning theory*, pp. 39–1. JMLR Workshop; Conference Proceedings.
- Bari BA, Moerke MJ, Jedema HP, Effinger DP, Cohen JY, Bradberry CW (2022) Reinforcement learning modeling reveals a reward-history-dependent strategy underlying reversal learning in squirrel monkeys. *Behavioral neuroscience*, **136**, 46. <https://doi.org/10.1037/bne0000492>
- Bartolo R, Averbach BB (2020) Prefrontal cortex predicts state switches during reversal learning. *Neuron*, **106**, 1044–1054. <https://doi.org/10.1016/j.neuron.2020.03.024>
- Berger-Tal O, Nathan J, Meron E, Saltz D (2014) The exploration-exploitation dilemma: A multidisciplinary framework. *PloS one*, **9**, e95693. <https://doi.org/10.1371/journal.pone.0095693>
- Bitterman ME (1975) The comparative analysis of learning: Are the laws of learning the same in all animals? *Science*, **188**, 699–709. <https://doi.org/10.1126/science.188.4189.699>
- Blaisdell A, Seitz B, Rowney C, Folsom M, MacPherson M, Deffner D, Logan CJ (2021b) Do the more flexible individuals rely more on causal cognition? Observation versus intervention in causal inference in great-tailed grackles. *Peer Community Journal*, **1**. <https://doi.org/10.24072/pcjournal.44>
- Blaisdell A, Seitz B, Rowney C, Folsom M, MacPherson M, Deffner D, Logan CJ (2021a) Do the more flexible individuals rely more on causal cognition? Observation versus intervention in causal inference in great-tailed grackles (version 5 of this preprint has been peer reviewed and recommended by peer community in ecology [<https://doi.org/10.24072/pci.ecology.100076>]). <https://doi.org/10.31234/osf.io/z4p6s>
- Bond AB, Kamil AC, Balda RP (2007) Serial reversal learning and the evolution of behavioral flexibility in three species of north american corvids (gymnorhinus cyanocephalus, nucifraga columbiana, aphelocoma californica). *Journal of Comparative Psychology*, **121**, 372. <https://doi.org/10.1037/0735-7036.121.4.372>
- Breen AJ, Deffner D (2023) Leading an urban invasion: Risk-sensitive learning is a winning strategy. *eLife*, **12**, RP89315. <https://doi.org/10.1101/2023.03.19.533319>
- Camerer C, Hua Ho T (1999) Experience-weighted attraction learning in normal form games. *Econometrica*, **67**, 827–874. <https://doi.org/10.1111/1468-0262.00054>
- Cauchoux M, Hermer E, Chaine A, Morand-Ferron J (2017) Cognition in the field: Comparison of reversal learning performance in captive and wild passerines. *Scientific reports*, **7**, 12945. <https://doi.org/10.1038/s41598-017-13179-5>
- Chen CS, Knep E, Han A, Ebitz RB, Grissom NM (2021) Sex differences in learning from exploration. *Elife*, **10**, e69748. <https://doi.org/10.7554/elife.69748>
- Chow PKY, Lea SE, Leaver LA (2016) How practice makes perfect: The role of persistence, flexibility and learning in problem-solving efficiency. *Animal behaviour*, **112**, 273–283. <https://doi.org/10.1016/j.anbehav.2015.11.014>
- Chow PK, Leaver LA, Wang M, Lea SE (2015) Serial reversal learning in gray squirrels: Learning efficiency as a function of learning and change of tactics. *Journal of Experimental Psychology: Animal Learning and Cognition*, **41**, 343. <https://doi.org/10.1037/xan0000072>
- Coulon A (2023) An experiment to improve our understanding of the link between behavioral flexibility and innovativeness. *Peer Community in Ecology*, **1**, 100407. <https://doi.org/10.24072/pci.ecology.100407>
- Danwitz L, Mathar D, Smith E, Tuzsus D, Peters J (2022) Parameter and model recovery of reinforcement learning models for restless bandit problems. *Computational Brain & Behavior*, **5**, 547–563. <https://doi.org/10.1007/s42113-022-00139-0>
- Daw ND, O’doherly JP, Dayan P, Seymour B, Dolan RJ (2006) Cortical substrates for exploratory decisions in humans. *Nature*, **441**, 876–879. <https://doi.org/10.1038/nature04766>
- Deffner D, Kleinow V, McElreath R (2020) Dynamic social learning in temporally and spatially variable environments. *Royal Society open science*, **7**, 200734. <https://doi.org/10.1098/rsos.200734>
- Degrande R, Cornilleau F, Lansade L, Jardat P, Colson V, Calandreau L (2022) Domestic hens succeed at serial reversal learning and perceptual concept generalisation using a new automated touchscreen device. *animal*, **16**, 100607. <https://doi.org/10.1016/j.animal.2022.100607>
- Donaldson-Matasci MC, Bergstrom CT, Lachmann M (2013) When unreliable cues are good enough. *The American Naturalist*, **182**, 313–327.
- Dufort RH, Guttman N, Kimble GA (1954) One-trial discrimination reversal in the white rat. *Journal of Comparative and Physiological Psychology*, **47**, 248. <https://doi.org/10.1037/h0057856>
- Dunlap AS, Stephens DW (2009) Components of change in the evolution of learning and unlearned preference. *Proceedings of the Royal Society B: Biological Sciences*, **276**, 3201–3208. <https://doi.org/10.1098/rspb.2009.0602>

- Erdsack N, Dehnhardt G, Hanke FD (2022) Serial visual reversal learning in harbor seals (*phoca vitulina*). *Animal Cognition*, **25**, 1183–1193. <https://doi.org/10.1007/s10071-022-01653-1>
- Federspiel IG, Garland A, Guez D, Bugnyar T, Healy SD, Güntürkün O, Griffin AS (2017) Adjusting foraging strategies: A comparison of rural and urban common mynas (*acridotheres tristis*). *Animal cognition*, **20**, 65–74. <https://doi.org/10.1007/s10071-016-1045-7>
- Frömer R, Nassar M (2023) Belief updates, learning and adaptive decision making. <https://doi.org/10.31234/osf.io/qndba>
- Gallistel CR, Fairhurst S, Balsam P (2004) The learning curve: Implications of a quantitative analysis. *Proceedings of the National Academy of Sciences*, **101**, 13124–13131. <https://doi.org/10.1073/pnas.0404965101>
- Gelman A, Rubin DB (1995) Avoiding model selection in bayesian social research. *Sociological methodology*, **25**, 165–173. <https://doi.org/10.2307/271064>
- Gershman SJ (2018) Deconstructing the human algorithms for exploration. *Cognition*, **173**, 34–42. <https://doi.org/10.1016/j.cognition.2017.12.014>
- Griffin AS, Guez D, Lermite F, Patience M (2013) Tracking changing environments: Innovators are fast, but not flexible learners. *PloS one*, **8**, e84907. <https://doi.org/10.1371/journal.pone.0084907>
- Izquierdo A, Brigman JL, Radke AK, Rudebeck PH, Holmes A (2017) The neural basis of reversal learning: An updated perspective. *Neuroscience*, **345**, 12–26. <https://doi.org/10.1016/j.neuroscience.2016.03.021>
- Komisichke B, Giurfa M, Lachnit H, Malun D (2002) Successive olfactory reversal learning in honeybees. *Learning & memory*, **9**, 122–129. <https://doi.org/10.1101/lm.44602>
- Kramer DL, Weary DM (1991) Exploration versus exploitation: A field study of time allocation to environmental tracking by foraging chipmunks. *Animal Behaviour*, **41**, 443–449. [https://doi.org/10.1016/s0003-3472\(05\)80846-2](https://doi.org/10.1016/s0003-3472(05)80846-2)
- Lea SE, Chow PK, Leaver LA, McLaren IP (2020) Behavioral flexibility: A review, a model, and some exploratory tests. *Learning & Behavior*, **48**, 173–187. <https://doi.org/10.3758/s13420-020-00421-w>
- Leimar O, Quiñones AE, Bshary R (2024) Flexible learning in complex worlds. *Behavioral Ecology*, **35**, arad109. <https://doi.org/10.1093/beheco/arad109>
- Liu Y, Day LB, Summers K, Burmeister SS (2016) Learning to learn: Advanced behavioural flexibility in a poison frog. *Animal Behaviour*, **111**, 167–172. <https://doi.org/10.1016/j.anbehav.2015.10.018>
- Logan CJ (2016) Behavioral flexibility in an invasive bird is independent of other behaviors. *PeerJ*, **4**, e2215. <https://doi.org/10.7717/peerj.2215>
- Logan C, Lukas D, Blaisdell A, Johnson-Ulrich Z, MacPherson M, Seitz B, Sevchik A, McCune K (2023a) Behavioral flexibility is manipulable and it improves flexibility and innovativeness in a new context. *Peer Community Journal*, **3**. <https://doi.org/10.24072/pcjournal.284>
- Logan C, Lukas D, Blaisdell A, Johnson-Ulrich Z, MacPherson M, Seitz B, Sevchik A, McCune K (2023b) Data: Behavioral flexibility is manipulable and it improves flexibility and problem solving in a new context. *Knowledge Network for Biocomplexity*, **Data package**. <https://doi.org/10.5063/F1BR8QNC>
- Logan C, Lukas D, Geng X, LeGrande-Rolls C, Marfori Z, MacPherson M, Rowney C, Smith C, McCune K (2024) Behavioral flexibility is related to foraging, but not social or habitat use behaviors in a species that is rapidly expanding its range. *EcoEvoRxiv*. <https://doi.org/10.32942/X2T036>
- Logan CJ, McCune KB, LeGrande-Rolls C, Marfori Z, Hubbard J, Lukas D (2023c) Implementing a rapid geographic range expansion - the role of behavior changes. *Peer Community Journal*. <https://doi.org/10.24072/pcjournal.320>
- Logan CJ, McCune K, MacPherson M, Johnson-Ulrich Z, Rowney C, Seitz B, Blaisdell A, Deffner D, Wascher C (2021) Are the more flexible great-tailed grackles also better at behavioral inhibition? *PsyArXiv*. <https://doi.org/10.31234/osf.io/vpc39>
- Logan CJ, Shaw R, Lukas D, McCune KB (2022) How to succeed in human modified environments. *In principle acceptance by PCI Ecology of the version on 8 Sep 2022*. <https://doi.org/https://doi.org/10.17605/OSF.IO/346AF>
- Mackintosh N, McGonigle B, Holgate V (1968) Factors underlying improvement in serial reversal learning. *Canadian Journal of Psychology/Revue canadienne de psychologie*, **22**, 85. <https://doi.org/10.1037/h0082753>
- McCune K, Blaisdell A, Johnson-Ulrich Z, Sevchik A, Lukas D, MacPherson M, Seitz B, Logan C (2023) Repeatability of performance within and across contexts measuring behavioral flexibility. *PeerJ*. <https://doi.org/10.7717/peerj.15773>
- McElreath R (2020) *Statistical rethinking: A bayesian course with examples in r and stan*. Chapman; Hall/CRC, Boca Raton, FL. <https://doi.org/10.1201/9780429029608>
- Metha JA, Brian ML, Oberrauch S, Barnes SA, Featherby TJ, Bossaerts P, Murawski C, Hoyer D, Jacobson

- LH (2020) Separating probability and reversal learning in a novel probabilistic reversal learning task for mice. *Frontiers in behavioral neuroscience*, **13**, 270.
- Mikhalevich I, Powell R, Logan C (2017) Is behavioural flexibility evidence of cognitive complexity? How evolution can inform comparative cognition. *Interface Focus*, **7**, 20160121. <https://doi.org/10.1098/rsfs.2016.0121>
- Minh Le N, Yildirim M, Wang Y, Sugihara H, Jazayeri M, Sur M (2023) Mixtures of strategies underlie rodent behavior during reversal learning. *PLOS Computational Biology*, **19**, e1011430. <https://doi.org/10.1371/journal.pcbi.1011430>
- Neftci EO, Averbach BB (2019) Reinforcement learning in artificial and biological systems. *Nature Machine Intelligence*, **1**, 133–143.
- R Core Team (2023) *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Rescorla RA, Wagner AR (1972) A theory of pavlovian conditioning: Variations in the effectiveness of reinforcement and nonreinforcement. In: *Classical conditioning II: Current theory and research* (eds Black AH, Prosy WF), pp. 64–99. Appleton-Century-Crofts, New York.
- Shettleworth SJ (2010) *Cognition, evolution, and behavior*. Oxford university press.
- Sih A (2013) Understanding variation in behavioural responses to human-induced rapid environmental change: A conceptual overview. *Animal Behaviour*, **85**, 1077–1088.
- Sol D, Timmermans S, Lefebvre L (2002) Behavioural flexibility and invasion success in birds. *Animal behaviour*, **63**, 495–502. <https://doi.org/10.1006/anbe.2001.1953>
- Spence KW (1936) The nature of discrimination learning in animals. *Psychological review*, **43**, 427. <https://doi.org/10.1037/h0056975>
- Stan Development Team (2023) *Stan modeling language users guide and reference manual, version 2.32.0*, <https://mc-stan.org/>.
- Strang CG, Sherry DF (2014) Serial reversal learning in bumblebees (*bombus impatiens*). *Animal Cognition*, **17**, 723–734. <https://doi.org/10.1007/s10071-013-0704-1>
- Tello-Ramos MC, Branch CL, Kozlovsky DY, Pitera AM, Pravosudov VV (2019) Spatial memory and cognitive flexibility trade-offs: To be or not to be flexible, that is the question. *Animal Behaviour*, **147**, 129–136. <https://doi.org/10.1016/j.anbehav.2018.02.019>
- Vehtari A, Gelman A, Simpson D, Carpenter B, Bürkner P-C (2021) Rank-normalization, folding, and localization: An improved rhat for assessing convergence of MCMC (with discussion). *Bayesian Analysis*. <https://doi.org/10.1214/20-BA1221>
- Warren J (1965a) Primate learning in comparative perspective. *Behavior of nonhuman primates*, **1**, 249–281. <https://doi.org/10.1016/B978-1-4832-2820-4.50014-7>
- Warren JM (1965b) The comparative psychology of learning. *Annual review of psychology*, **16**, 95–118. <https://doi.org/10.1146/annurev.ps.16.020165.000523>
- Warren J, Warren HB (1962) Reversal learning by horse and raccoon. *The Journal of Genetic Psychology*, **100**, 215–220. <https://doi.org/10.1080/00221325.1962.10533590>
- Woo JH, Aguirre CG, Bari BA, Tsutsui K-I, Grabenhorst F, Cohen JY, Schultz W, Izquierdo A, Soltani A (2023) Mechanisms of adjustments to different types of uncertainty in the reward environment across mice and monkeys. *Cognitive, Affective, & Behavioral Neuroscience*, 1–20. <https://doi.org/10.1101/2022.10.01.510477>
- Wright TF, Eberhard JR, Hobson EA, Avery ML, Russello MA (2010) Behavioral flexibility and species invasions: The adaptive flexibility hypothesis. *Ethology Ecology & Evolution*, **22**, 393–404. <https://doi.org/10.1080/03949370.2010.505580>