

## Highlights

### **Learning Rates Are Not All the Same: The Interpretation of Computational Model Parameters Depends on the Context**

Maria K. Eckstein, Sarah L. Master, Liyu Xia, Ronald E. Dahl, Linda Wilbrecht, Anne G.E. Collins

- Efforts in computational cognitive modeling often assume that Reinforcement Learning (RL) modeling parameters will generalize between studies and models, but this is not well established.
- We empirically investigate whether RL parameters generalize between three tasks and models, using a large developmental dataset and a within-participant design.
- We find that RL decision noise/exploration parameters generalize fairly well, but RL learning rates do not.
- Our data support previous conclusions that decision noise/exploration decreases during development (ages 8-17), but suggests that claims about learning rate development cannot be generalized.

# Learning Rates Are Not All the Same: The Interpretation of Computational Model Parameters Depends on the Context

Maria K. Eckstein<sup>a</sup>, Sarah L. Master<sup>a,b</sup>, Liyu Xia<sup>a,c</sup>, Ronald E. Dahl<sup>a,d</sup>,  
Linda Wilbrecht<sup>a,e</sup>, Anne G.E. Collins<sup>a,e</sup>

<sup>a</sup>*Department of Psychology, UC Berkeley, 2121 Berkeley Way  
West, Berkeley, 94720, CA, USA*

<sup>b</sup>*Department of Psychology, New York University, 6 Washington Place, New  
York, 10003, NY, USA*

<sup>c</sup>*Department of Mathematics, UC Berkeley, 970 Evans Hall, Berkeley, 94720, CA, USA*

<sup>d</sup>*Institute of Human Development, UC Berkeley, 2121 Berkeley Way  
West, Berkeley, 94720, CA, USA*

<sup>e</sup>*Helen Wills Neuroscience Institute, UC Berkeley, 175 Li Ka Shing  
Center, Berkeley, 94720, CA, USA*

---

## Abstract

Reinforcement Learning (RL) has revolutionized the cognitive and brain sciences, explaining behavior from simple conditioning to problem solving, across the life span, and anchored in brain function. However, discrepancies in results are increasingly apparent between studies, particularly in the developmental literature. To better understand these, we investigated to which extent parameters *generalize* between tasks and models, and capture specific and uniquely *interpretable* (neuro)cognitive processes. 291 participants aged 8-30 years completed three learning tasks in a single session, and were fitted using state-of-the-art RL models. RL decision noise/exploration parameters generalized well between tasks, decreasing between ages 8-17. Learning rates for negative feedback did not generalize, and learning rates for positive feedback showed intermediate generalizability, dependent on task similarity. These findings can explain discrepancies in the existing literature. Future research therefore needs to carefully consider task characteristics when relating findings across studies, and develop strategies to computationally model how context impacts behavior.

*Keywords:* Computational Modeling, Model Parameter, Reinforcement

## Learning, Development, Generalizability, Interpretability

---

### 1 **1. Introduction**

2 In recent decades, the cognitive neurosciences have made breakthroughs  
3 in computational modeling, showing that reinforcement learning (RL) mod-  
4 els can explain foundational aspects of human behavior. RL does not only  
5 seem to underlie simple cognitive processes such as stimulus-outcome and  
6 stimulus-response learning [1, 2, 3], but also complex ones, including goal-  
7 directed, temporally-extended behavior [4, 5], meta-learning [6], and abstract  
8 problem solving that requires hierarchical thinking [7, 8, 9, 10]. Underlining  
9 their centrality in the study of human cognition, RL models have been ap-  
10 plied across the lifespan [11, 12, 13], and in healthy participants as well those  
11 experiencing psychiatric illnesses [14, 15, 16, 17, 18]. RL models are of partic-  
12 ular interest because they also capture brain function: A specialized network  
13 of brain regions, including the basal ganglia and prefrontal cortex, implement  
14 computations that mirror specific components of RL algorithms, including  
15 action values and reward prediction errors [19, 20, 21, 22, 23, 24, 25]. In sum,  
16 explaining behaviors from simple conditioning to complex problem solving,  
17 adequate for diverse human populations, based on a compelling theoretical  
18 foundation [26], and with strong ties to brain function, RL has experienced  
19 a surge in published studies since its inception [27], and emerged as a pow-  
20 erful and potentially unifying modeling framework for cognitive and neural  
21 processing.

22 Despite their increasing popularity, however, not enough attention has  
23 been paid to what exactly RL models and model variables (e.g., model pa-  
24 rameters) measure, and our current assumptions might be imprecise, po-  
25 tentially slowing further progress. Our recent opinion paper develops this  
26 argument in depth [28]. In brief, computational modeling condenses be-  
27 havioral datasets into a model and a small number of free model parameters  
28 [11, 27, 29, 30, 31, 32]. We as researchers often assume that these models and  
29 parameters expose mental and/or neural processes, and have the ability to  
30 dissect them into specific, unique components (e.g., value updating and deci-  
31 sion making), thereby measuring participants' inherent characteristics (e.g.,  
32 individual learning rates). However, we argue in this paper that these as-  
33 sumptions might be too optimistic and that a careful empirical investigation  
34 is required to assess their validity.

35 We focus on two major aspects, which are adopted widely in computa-  
36 tional modeling [28]: *generalizability* and *interpretability*. We define a model  
37 variable (e.g., fitted parameter) as *generalizable* if it is consistent across uses,  
38 such that a person would be characterized with the same values independent  
39 of the specific model or task used to estimate the variable. Generalizability is  
40 a consequence of the assumption that parameters are intrinsic to participants  
41 (e.g., a person with a high learning rate) rather than task dependent. We  
42 further define a model parameter as *interpretable* if it isolates specific and  
43 unique elements of cognition, which are often assumed to be implemented in  
44 separable neural substrates: Decomposing behavior into model parameters  
45 is seen as a way of *carving cognition at its joints*.

46 Assumptions about generalizability and interpretability are rarely stated  
47 explicitly, but underlie conclusions across the fields of computational psy-  
48 chology and neuroscience, and often implicitly guide research efforts. As-  
49 sumptions of generalizability, for example, inspired many to identify the  
50 inherent, task-independent settings of parameters in humans (e.g., empiri-  
51 cal parameter distributions [33]; relationships between negative and positive  
52 learning rates [34]), to characterize the age development of parameters in a  
53 task-independent way [11, 12, 13, 35], and to compare parameters between  
54 studies in review articles [14, 15, 16, 19, 20, 21, 22, 23, 25], meta-analyses  
55 [24, 36, 37], and discussion sections of empirical papers: When model vari-  
56 ables are compared between different types of studies, there is an implicit  
57 assumption of generalization. Relying on interpretability, model variables  
58 have been expected to be associated with specific neural substrates (e.g.,  
59 reward prediction errors and dopamine function [38]), to expose the core of  
60 what differentiates participants with psychiatric conditions from healthy ones  
61 (e.g., working-memory parameter differences in schizophrenia [39]), and gener-  
62 ally, to capture processes that are particularly “theoretically meaningful”  
63 [14].

64 However, inconsistencies in empirical results are emerging across the de-  
65 velopmental [13, 40, 41, 42], clinical [15, 16, 17, 18], cognitive, and neuroscien-  
66 tific literature [24, 36, 37, 43], potentially suggesting a lack of generalizability  
67 and/or interpretability, which is also in accordance with different theoretic-  
68 al considerations [27, 28, 44, 45, 46, 47]. Nevertheless, the degree to which  
69 parameters generalize between tasks and are interpretable has not been in-  
70 vestigated empirically yet (but see [48] for work on parameter reliability).

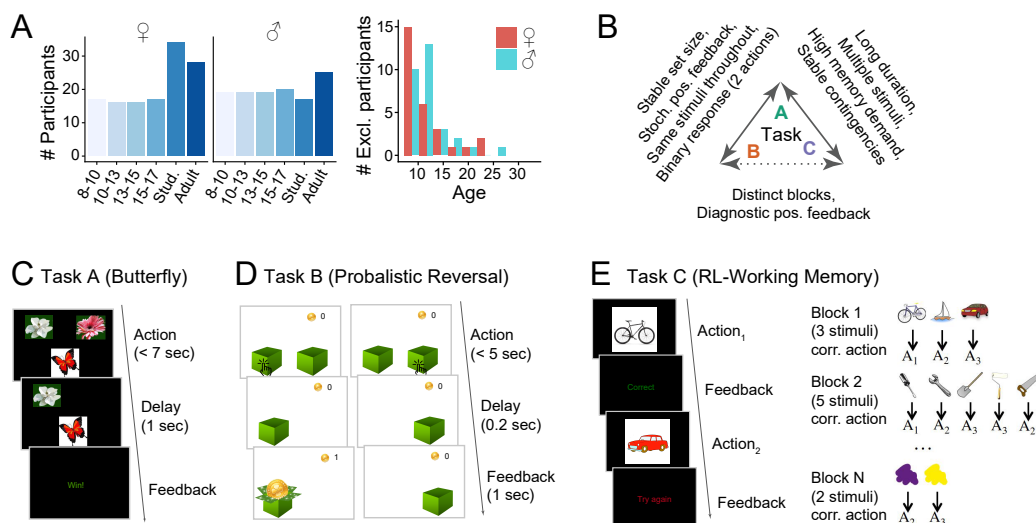
71 This was the goal of the current project. We compared the RL parameters  
72 fit to the same individuals across different learning tasks in a single study.

73 We used a developmental dataset (291 participants, ages 8-30 years), which  
74 allowed us to obtain a large spread of individual differences and address  
75 outstanding discrepancies in the developmental psychology literature [13].  
76 The three learning tasks varied on several common dimensions, including  
77 feedback stochasticity, task volatility, and memory demands (Fig. 1B), and  
78 have previously been used to study RL processes [35, 49, 50]. However, like  
79 many tasks in the literature, these tasks likely also engaged other cognitive  
80 processes, such as working memory and reasoning. The within-participant  
81 design allowed us to test directly whether the same participants showed the  
82 same parameters across tasks (generalizability), and the combination of mul-  
83 tiple tasks shed light on which cognitive processes parameters captured in  
84 each task (interpretability). We extensively compared and validated the RL  
85 models of each task [27, 30, 51], and previously reported the developmental  
86 results separately [35, 49, 50].

87 We found that the RL parameters that reflect decision noise or exploration  
88 (inverse decision temperature  $\beta$ , undirected noise  $\epsilon$ ; for model details, see  
89 section 4.5) were most consistent within individuals across tasks, suggesting  
90 that these parameters were most generalizable. Decision noise/exploration  
91 parameters also showed a consistent developmental pattern across subjects,  
92 declining from age 8-17. RL learning rate parameters ( $\alpha_+$ ,  $\alpha_-$ ), however,  
93 were largely inconsistent within individuals across tasks, showing that they  
94 did not generalize. Capturing different variance, they likely also reflected  
95 different cognitive processes across tasks. Both of these patterns are con-  
96 sistent with patterns that have started to emerge in the existing literature  
97 [13]. Behavioral analyses indicated that task differences, and the associated  
98 differences in optimal behavior, might underlie these observed parameter dis-  
99 crepancies. These results suggest that past computational findings are not  
100 as generalizable as often assumed, and that future research needs to address  
101 the reasons of the observed discrepancies to move the field forward.

## 102 2. Results

103 The next section gives a brief overview of the experimental tasks and  
104 computational models, before tackling parameter generalizability (section  
105 2.1) and parameter interpretability (section 2.2). Task details are provided in  
106 Fig. 1C-E and section 4.4, and computational models and parameter fitting  
107 in section 4.5, as well as the original publications [35, 49, 50].



**Figure 1: Overview of the experimental paradigm.** (A) Participant sample. Left: Number of participants in each age group, broken up by sex (self-reported). Age groups were determined by within-sex age quartiles for participants between 8-17 years. The adult sample is broken up by recruitment type (“Stud.”: University undergraduates, receiving course credit for participation. “Adult”: Adults recruited from the community using the same methods as the developing participants). Right: Number of participants who participated in the study and whose data were excluded because they failed to reach the performance criterion in at least one task. (B) Pairwise similarities in task design between tasks A, B, and C. Similarities between each pair of tasks are shown above the connecting arrows. Only features are shown that differentiate two tasks from the third. E.g., noting “Stable set size” on the edge between tasks A and B implies that set size was not stable in task C. Task A shared more similarities with tasks B and C than they shared with each other. (C) Procedure of task A (“Butterfly task”). Participants saw one of four butterflies on each trial and selected one of two flowers in response, via button press on a game controller. Each butterfly had a stable preference for one flower throughout the task, but rewards were delivered stochastically (70% for correct responses, 30% for incorrect). For details, see section 4.4 and the original publication [50]. (D) Procedure of task B (“Probabilistic switching”). Participants saw two boxes on each trial and selected one with the goal of finding gold coins. At each point in time, one box was correct and had a high (75%) probability of delivering a coin, whereas the other was incorrect (0%). At unpredictable intervals, the correct box switched sides. For details, see section 4.4 and [49]. (E) Procedure of task C (“Reinforcement learning-working memory”). Participants saw one stimulus on each trial and selected one of three responses. All correct responses and no incorrect responses were rewarded. Stimuli were presented in blocks containing 2-5 different stimuli. The number of stimuli in a block is called set size. The task was designed to disentangle set-size sensitive working memory processes from set-size insensitive RL processes. For details, see section 4.4 and [35].

108 Depending on the task, RL models contained different parameters, re-  
109 flecting existing differences in the literature. Task A required participants  
110 to learn the correct associations between each of four stimuli (butterflies)  
111 and two responses (flowers), through probabilistic feedback (Fig. 1C). The  
112 best-fitting model contained three free parameters: learning rate from posi-  
113 tive outcomes  $\alpha_+$ , inverse decision temperature  $\beta$ , and Forgetting  $F$ ; and one  
114 fixed parameter: learning rate from negative outcomes  $\alpha_- = 0$  [50]. Task B  
115 required participants to adapt to unexpected switches in the action-outcome  
116 contingencies of a simple bandit task (only one of two boxes contained a  
117 gold coin at any time), based on semi-probabilistic feedback (Fig. 1D). The  
118 best-fitting RL model contained four free parameters:  $\alpha_+$ ,  $\alpha_-$ ,  $\beta$ , and choice  
119 persistence  $p$  [49]. Task C required learning of stimulus-response associations  
120 like task A, but over several task blocks with varying numbers of stimuli,  
121 and provided deterministic feedback (Fig. 1E). The best model for this task  
122 combined RL and working-memory processes, containing RL parameters  $\alpha_+$   
123 and  $\alpha_-$ ; working-memory parameters capacity  $K$ , Forgetting  $F$ , and noise  $\epsilon$ ;  
124 and mixture parameter  $\rho$ , which determined the relative weights of RL and  
125 working memory [35, 52].

126 To ensure that potential parameter discrepancies in this study were not  
127 due to a lack of modeling quality, we employed rigorous model fitting, com-  
128 parison, and validation [27, 29, 30, 51]: For each task, we compared a large  
129 number of competing models, based on different parameterizations and cog-  
130 nitive mechanisms, and selected the best one based on quantitative model  
131 comparison scores, models' ability to reproduce participants' behavior in sim-  
132 ulation, and other criteria of model fit (e.g., interpretability) [44, 45]. We  
133 also used hierarchical Bayesian methods for model fitting and comparison  
134 when possible to obtain most accurate parameter estimates [51]. Individual  
135 publications provide further details [35, 49, 50].

### 136 *2.1. Part I: Parameter Generalizability*

137 To investigate parameter generalizability, we assessed whether partici-  
138 pants showed similar parameter values across tasks, and whether different  
139 tasks showed the same parameter age trajectories. These within-participant  
140 comparisons are crucial to determine whether discrepancies in the previ-  
141 ous literature were caused by methodological differences (e.g., differences in  
142 participant samples, testing procedures, modeling quality, research labs), or  
143 could arise from mere differences in task characteristics and computational  
144 models, as we hypothesized.



145 *2.1.1. Differences in Absolute Parameter Values*

146 We first asked whether tasks led to different absolute parameter val-  
147 ues (Fig. 2A), using repeated-measures analyses of variance (ANOVAs).  
148 When ANOVAs showed significant task effects, we followed up with pair-  
149 wise, repeated-measures t-tests, using the Bonferroni correction.

150 Learning rates  $\alpha_+$  and  $\alpha_-$  occupied largely distinct ranges across tasks:  
151 Values were very low in tasks C ( $\alpha_+$  mean: 0.07, sd: 0.18;  $\alpha_-$  mean: 0.03, sd:  
152 0.13), intermediate in task A ( $\alpha_+$  mean: 0.22, sd: 0.09;  $\alpha_-$  was fixed at 0),  
153 and fairly high in task B ( $\alpha_+$  mean: 0.77, sd: 0.11;  $\alpha_-$  mean: 0.62, sd: 0.14;  
154 for statistical comparisons, see Table 1). Decision noise was high in task B ( $\frac{1}{\beta}$   
155 mean: 0.33, sd: 0.15), but low in tasks A ( $\frac{1}{\beta}$  mean: 0.095, sd: 0.0087) and C  
156 ( $\epsilon$  mean: 0.025, sd: 0.032; statistics in Table 1 ignore  $\epsilon$  because its absolute  
157 values were not comparable to  $\frac{1}{\beta}$  due to the different parameterization; see  
158 section 4.5). Forgetting was significantly higher in task C (mean: 0.19, sd:  
159 0.17) than A (mean: 0.056, sd: 0.028; task B was best fit without forgetting).

160 All ANOVAs revealed significant and large task effects, and all follow-up  
161 t-tests revealed significant and large pairwise differences (Table 1), showing  
162 that absolute parameter values differed substantially between tasks. This  
163 shows that the three tasks produced significantly different estimates of learn-  
164 ing rate, decision noise/exploration, and forgetting for the same participants  
165 (Fig. 2B). Interestingly, these parameter differences echoed differences in  
166 task demands: Learning rates and noise/exploration were highest in task B,  
167 where frequent switches required quick updating and high levels of explo-  
168 ration. Similarly, forgetting was highest in task C, which posed the largest  
169 demands on memory. Using regression models that controlled for age (instead  
170 of ANOVA) led to similar results (Table D.9).

171 *2.1.2. Relative Parameter Differences*

172 However, comparing absolute parameter values between tasks has short-  
173 comings: It ignores variance between participants, even though between-  
174 participant variance might be the more meaningful measure because it re-  
175 flects participants' relationships to each other. The simplest way to inves-  
176 tigate whether between-participant variance generalized between tasks is to  
177 test if individual variance in one task mirrors individual variance in another,  
178 using Spearman correlation (suppl. Fig. D.8). Indeed, both  $\alpha_+$  (suppl. Fig.  
179 D.8A) and noise/exploration parameters (suppl. Fig. D.8B) were signifi-  
180 cantly positively correlated between task A and tasks B and C, suggesting



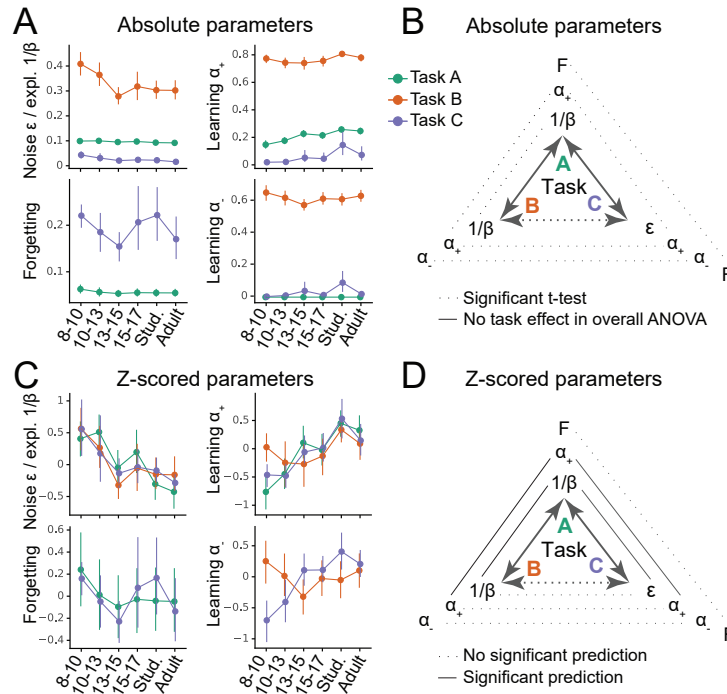


Figure 2: **Generalizability of absolute parameter values (A-B) and of parameter age trajectories / z-scored parameters (C-D) between tasks.** (A) Fitted parameters over participant age (quartile bins), for all three tasks (A: green; B: orange; C: blue). Parameter values differed significantly between tasks (for statistics, see Table 1). Dots indicate means, error bars specify the confidence level (0-1) for interval estimation of the population mean. (B) Summary of the main results of part (A), visualizing Table 1. Double-sided arrows are replicated from Fig. 1B and show task similarity. Lines show test statistics for absolute parameter values. Dotted lines indicate significant task differences in Bonferroni-corrected pairwise t-tests, which were conducted after observing significant task effects in corresponding ANOVAs. All t-tests were significant, indicating that absolute parameter values differed significantly for each pair of tasks. (C) Parameter age trajectories, i.e., within-task z-scored parameters over participant age bins. Age trajectories can potentially reveal similarities that are obscured by differences in means or variances when assessing absolute parameter values. (D) Summary of the main results of part (C), visualizing Table 4. When parameters in two tasks are connected with a full line, the parameter can be predicted significantly in one task from the other. When parameters are connected with a dotted line, the prediction is not significant. In contrast to absolute parameter values, age trajectories were predictive in several cases, especially for tasks with more similarities (A and B; A and C), compared to tasks with fewer (B and C).

181 that variance between participants generalized better than absolute values.  
182 However, significant correlations were lacking between tasks B and C. This  
183 suggests that  $\alpha_+$  and noise/exploration generalized from and to task A, but  
184 they did not generalize between tasks B and C, mirroring task similarities  
185 (Fig. 1B; also see section 2.2.1; Fig. D.9 shows the correlations between  
186 all pairs of features in the dataset.) Note that noise parameters generalized  
187 between task A and C despite differences in parameterization ( $\epsilon$  vs.  $1/\beta$ ),  
188 showing robustness in the characterization of choice stochasticity (suppl. Fig.  
189 D.8B).

### 190 2.1.3. Parameter Age Trajectories

191 However, this correlation analysis is limited in its failure to take into  
192 account age effects, a known source of variance, such that apparent task  
193 similarities could be driven by a shared dependence on age rather than age-  
194 independent underlying similarities. To address this, we next analyzed pa-  
195 rameters' age trajectories, which allowed us to abstract away potentially  
196 arbitrary differences (e.g., different parameter means and variances across  
197 tasks), while conserving potentially meaningful structure in the dataset (i.e.,  
198 participants' parameter values relative to each other).

199 We obtained age trajectories by z-scoring each parameter within each  
200 task (Fig. 2C). To test for differences in age trajectories, we used mixed-  
201 effects regression to predict parameters of all tasks from two age predictors  
202 (age and squared age) and task (A, B, or C). When this model fit better  
203 than the corresponding model without task, task characteristics affected age  
204 trajectories, and we added post-hoc models for each pair of tasks.

205 For  $\alpha_-$ , the task-based regression model showed a significantly better fit,  
206 revealing significant task differences (Table 2). Indeed,  $\alpha_-$  showed funda-  
207 mentally different age trajectories in task B compared to C (in task A,  $\alpha_-$   
208 was fixed): In task B,  $\alpha_-$  decreased linearly, modulated by a U-shaped cur-  
209 vature (linear effect of age:  $\beta = -0.11$ ,  $p < 0.001$ ; quadratic:  $\beta = 0.003$ ,  
210  $p < 0.001$ ), but in task C, it increased linearly, modulated by an inverse-U  
211 curvature (linear:  $\beta = 0.32$ ,  $p < 0.001$ ; quadratic:  $\beta = -0.07$ ,  $p < 0.001$ ;  
212 Fig. 2C). These differences were reflected in the significant interaction terms  
213 of the grand regression model (Table ??).

214 For  $\alpha_+$ , adding task as a predictor did not improve model fit, suggesting  
215 that age trajectories did not differ (Table 2). Indeed, age trajectories were  
216 qualitatively similar between tasks, showing linear increases that tapered off  
217 with age (linear increase: task A:  $\beta = 0.33$ ,  $p < 0.001$ ; task B:  $\beta = 0.052$ ,

218  $p < 0.001$ ; task C:  $\beta = 0.28$ ,  $p < 0.001$ ; quadratic modulation: task A:  
219  $\beta = -0.007$ ,  $p < 0.001$ ; task B:  $\beta = -0.001$ ,  $p < 0.001$ ; task C:  $\beta = -0.006$ ,  
220  $p < 0.001$ ).

221 For noise/exploration and Forgetting parameters, age trajectories did not  
222 differ either (Table 2). For decision noise/exploration, the grand regression  
223 model revealed a linear decrease and tapering off in older participants that  
224 was consistent across all tasks (Fig. 2C; Table ??), in accordance with previ-  
225 ous findings [13]. For Forgetting, the grand model did not reveal consistent  
226 age effects (Fig. 2C; Table ??).

227 In summary, when assessing absolute parameter values (Fig. 2A, 2B), dif-  
228 ferences in scale obscured existing similarities in age trajectories for noise/exploration  
229 parameters and  $\alpha_+$  (Fig. 2C). For  $\alpha_-$ , on the other hand, differences existed  
230 both in terms of scale (Fig. 2A, 2B) and age trajectories (Fig. 2C). As sug-  
231 gested by the correlation analysis, patterns of generalization differed between  
232 pairs of tasks, such that more generalization was present between tasks that  
233 were more similar in terms of task characteristics (A and B; A and C; not B  
234 and C).

#### 235 2.1.4. Predicting Age Trajectories

236 So far, we have assessed parameter differences to reveal parameters that  
237 do not generalize across tasks. However, the absence of differences only pro-  
238 vides indirect evidence *for* generalization. We therefore next assessed how  
239 closely parameters were related, using linear regression to predict partici-  
240 pants' parameters in one task from the values of the same parameter in a  
241 different task. We controlled for age by including age and squared age as  
242 predictors to ensure that the prediction was driven by parameter similarities  
243 beyond age.

244 For both  $\alpha_+$  and noise/exploration parameters, task A predicted tasks B  
245 and C, and tasks B and C predicted task A, but tasks B and C did not pre-  
246 dict each other (Table 4; Fig. 2D), confirming that  $\alpha_+$  and noise/exploration  
247 generalized from and to task A, but not between tasks B and C, mirroring  
248 task similarities (Fig. 1B; also see section 2.2.1). For  $\alpha_-$ , tasks B and C  
249 showed a marginally significant *negative* relationship (Table 4), suggesting  
250 that predicting  $\alpha_-$  in one task from the other would lead to inverse predic-  
251 tions. Indeed, we previously reported a U-shaped trajectory of  $\alpha_-$  in task  
252 B with minimum in 13-to-15-year-olds [49], but a consistent increase up to  
253 early adulthood in task C [50], revealing striking qualitative differences in  
254 the estimation of  $\alpha_-$  when using these two tasks. For Forgetting, tasks A

Table 1: Statistics of ANOVAs predicting raw parameter values from task (A, B, C). When an ANOVA showed a significant task effect, post-hoc, Bonferroni-corrected t-tests were added. \*  $p < .05$ ; \*\*  $p < .01$ , \*\*\*  $p < .001$ .

Parameter	Model	Tasks	F / t	df	$p$	sig.
$\frac{1}{\beta}$	ANOVA	A, B	830	1	$p < 0.001$	***
	t-test	A vs B	25	246	$p < 0.001$	***
$\alpha_+$	ANOVA	A, B, C	2.018	2	$p < 0.001$	***
	t-test	A vs B	66	246	$p < 0.001$	***
	t-test	A vs C	12	246	$p < 0.001$	***
	t-test	B vs C	51	246	$p < 0.001$	***
$\alpha_-$	ANOVA	B, C	2.357	1	$p < 0.001$	***
	t-test	B vs C	49	246	$p < 0.001$	***
Forgetting	ANOVA	A, C	161	1	$p < 0.001$	***
	t-test	A vs C	49	246	$p < 0.001$	***

Table 2: Assessing the existence of age effects on parameter trajectories: Model fits of regression models predicting parameter age trajectories, comparing the added value of including (“AIC with task”) versus excluding (“AIC without task”) task as a predictor. Differences in AIC scores were tested statistically using F-tests. The best (significantly smaller) AIC scores are highlighted in bold, and their coefficients are shown in Table ??.

Parameter	AIC without task	AIC with task	F(df)	p	sig.
$\frac{1}{\beta}/\epsilon$	<b>2,044</b>	2,054	NA	NA	–
$\alpha_+$	<b>2,044</b>	2,042	$F(4, 245) = 2.34$	$p = 0.056$	–
$\alpha_-$	1,395	<b>1,373</b>	$F(2, 245) = 6.99$	$p = 0.0011$	**
Forgetting	<b>1,406</b>	1,411	NA	NA	–

Table 3: Statistical tests on age trajectories: mixed-effects regression models predicting z-scored parameter values from task (A, B, C), age, and squared age (months). When the task-less model fitted best, the coefficients of this model are shown, showing shared age trajectories (Table 2;  $\frac{1}{\beta}/\epsilon$ ,  $\alpha_+$ , Forgetting). When the age-based model fitted better, pairwise follow-up models are shown ( $\alpha_-$ ), showing task differences. P-values of follow-up models were corrected for multiple comparison using the Bonferroni correction. \*  $p < .05$ ; \*\*  $p < .01$ , \*\*\*  $p < .001$ .

Parameter	Tasks	Predictor	$\beta$	$p$ (Bonf.)	sig.
$\frac{1}{\beta}/\epsilon$	A, B, C	Intercept	1.86	< 0.001	***
		Age (linear)	-0.17	0.003	**
		Age (quadratic)	0.004	< 0.001	***
$\alpha_+$	A, B, C	Intercept	-2.10	< 0.001	***
		Age (linear)	0.20	< 0.001	***
		Age (quadratic)	-0.004	< 0.001	***
$\alpha_-$	B, C	Task (main effect)	4.15	< 0.001	***
		Task * linear age (interaction)	0.43	< 0.001	***
		Task * quadratic age (interaction)	-0.010	< 0.001	***
Forgetting	A, C	Intercept	0.37	0.44	
		Age (linear)	-0.034	0.53	
		Age (quadratic)	0.001	0.63	

255 and C were not predictive of each other (Table 4).

256 Importantly, these results (Fig. 2D) differ from the previous patterns  
 257 (Fig. 2C) for Forgetting parameters and  $\alpha_+$  in tasks B and C. This shows  
 258 that a lack of difference (Fig. 2C) does not imply successful prediction (Fig.  
 259 2D).

### 260 2.1.5. Summary Part I

261 In summary, Part I revealed that (1) different tasks led to different es-  
 262 timates of participants' exploration ( $\frac{1}{\beta}$ ), Forgetting ( $F$ ), and learning rates  
 263 ( $\alpha_+$ ,  $\alpha_-$ ), revealing a lack of generalization of absolute parameter values.  
 264 Intriguingly, absolute parameter values were stable within tasks (reflecting  
 265 task demands), but varied within participants. (2) In contrast to absolute  
 266 parameter values, age trajectories of noise/exploration parameters and learn-  
 267 ing rates  $\alpha_+$  were qualitatively similar between tasks, suggesting that pa-  
 268 rameter age trajectories generalized better than absolute values. The age  
 269 trajectories of learning rates  $\alpha_-$ , however, differed fundamentally between  
 270 tasks, highlighting that parameters in the same models can generalize dif-  
 271 ferently. (3) Assessing the parameters with task-consistent age trajectories,

272 noise/exploration decreased until early adulthood, in accordance with the  
273 literature [13], while learning rates  $\alpha_+$  increased. (4) Using between-task  
274 prediction as the strongest test of generalization, age trajectories of learning  
275 rates  $\alpha_-$  and Forgetting were not predictive, and noise/exploration param-  
276 eters and learning rates  $\alpha_+$  could only be predicted between similar tasks,  
277 suggesting that generalizability was generally weaker than expected, and  
278 might depend on task similarity.

## 279 *2.2. Part II: Parameter Interpretability*

280 Based on these insights, Part II of our investigations focused on param-  
281 eter interpretability, i.e., the concept that parameters capture specific and  
282 unique cognitive processes that are well delineated. We tested parameter in-  
283 terpretability by investigating the relations between different parameters in  
284 our dataset, assessing the specificity and distinctiveness of each parameter as  
285 well as the relations between parameters and observed patterns of behavior.

### 286 *2.2.1. The Main Axes of Variation*

287 To gain an understanding of what information was captured by each  
288 parameter, we employed a data-driven approach, identifying major axes of  
289 variance without specifying a priori hypotheses. We used PCA to identify the  
290 major axes in our dataset (composed of both behavioral features and model  
291 parameters). We then used these axes (principal components; PCs) to in-  
292 terpret model parameters. To understand the PCs themselves, we analyzed  
293 the weights of the behavioral features on each PC (Fig. 3). Detailed infor-  
294 mation is provided in sections 4.6 (PCA methods), Appendix C (behavioral  
295 features), and suppl. Fig. D.10 (additional PCA results).

296 We first examined PC1, the axis of largest variation (25.1% of explained  
297 variance; suppl. Fig. D.10A), to understand the main sources of individual  
298 differences in our dataset. Behaviors that indicated good task participation  
299 (e.g., high percentage of correct choices) loaded positively on PC1, whereas  
300 behaviors that indicated that participants were not on task loaded negatively  
301 (e.g., more missed trials, longer response times; Fig. 3A). PC1 comprised  
302 measures both in the narrow sense of maximizing task accuracy (e.g., per-  
303 centage correct choices, measures of task accuracy, win-stay choices), and in  
304 the wider sense of reflecting task engagement (e.g., number of missed trials,  
305 response times, response time variability). PC1 therefore captured a range  
306 of “good performance” indicators, reflecting general task engagement. PC1

307 increased significantly with age, consistent with participants' increasing per-  
308 formance (suppl. Fig. B.6B; age effects of subsequent PCs in suppl. Fig.  
309 D.10; suppl. Table D.8).

310 In all three tasks, noise/exploration loaded negatively on PC1 (Fig. 3A),  
311 showing that elevated decision stochasticity was associated with poorer per-  
312 formance in all tasks. Forgetting parameters also loaded negatively, support-  
313 ing a negative role for performance.  $\alpha_+$  showed positive loadings in all three  
314 tasks, suggesting that faster integration of positive feedback was associated  
315 with better performance. Intriguingly,  $\alpha_-$  loaded positively in task C, but  
316 negatively in task B, suggesting that performance increased when partici-  
317 pants integrated negative feedback faster in task C, but decreased when they  
318 did the same in task B. This distinction can be interpreted in terms of task  
319 demands: Negative feedback was diagnostic in task C, but non-diagnostic in  
320 task B (Fig. 1B), such that repeating choices after negative feedback ("Lose-  
321 stay" behavior) was hurtful in the former (negative loading on PC1 for task  
322 C), but can be beneficial in the latter (positive loading on PC1 for task B;  
323 Fig. 3A).

324 Having gained insight into parameters' roles for task engagement by an-  
325 alyzing PC1, we next turned to PC2 and PC3. To facilitate their interpreta-  
326 tion, we flipped the loadings of all PC2 and PC3 features that were negative  
327 on PC1, to make them interpretable with respect to task engagement (for  
328 methodological details, see section 4.6). This pre-processing revealed that  
329 PC2 and PC3 encoded task contrasts: PC2 contrasted task B to task C  
330 (loadings on corresponding features were positive / negative / near-zero for  
331 tasks B / C / A; Fig. 3B). PC3 contrasted task A to both B and C (load-  
332 ings on corresponding features were positive / negative for task A / tasks  
333 B and C; Fig. 3C; missed trials and response times did not show task con-  
334 trasts, suggesting that these features did not differentiate between tasks).  
335 The ordering of PC2 and PC3 shows that participants' behavior differed  
336 more between tasks B and C (PC2: 8.9% explained variance) than between  
337 B or C and A (PC3: 6.2%; suppl. Fig. D.10), in accordance with descriptive  
338 task characteristics (Fig. 1B). This shows that after task engagement, the  
339 main variation in our dataset arose from task differences.

340 Intriguingly, noise/exploration parameters,  $\alpha_+$ , and  $\alpha_-$  reproduced the  
341 task contrasts of PC2 and PC3, showing positive or negative loadings based  
342 on the task in which they were measured (Fig. 3B, 3C). This means that these  
343 parameters differed sufficiently between tasks to be discriminable (as opposed  
344 to, e.g., response times and numbers of missed trials, which did not show



345 task contrasts, suggesting that they were not discriminable between tasks).  
346 Each parameter therefore captured enough task-specific variance to make it  
347 possible to be identified with the correct task. This degree of differentiability  
348 would not be expected if parameters captured the same processes in each  
349 task, in which case they would capture the same variance and not show  
350 task differences. Taken together, PC2 and PC3 confirmed that each of these  
351 parameters captured task-unique processes.

352 Taken together, the PCA revealed that (1) the main axes of variation in  
353 the dataset were task engagement (PC1) and task differences (PC2-PC3).  
354 (2) Noise/exploration, Forgetting,  $\alpha_+$ , and  $\alpha_-$  all were related to task en-  
355 gagement (PC1). Whereas the relation was consistent between tasks for the  
356 former three, it was task-dependent for  $\alpha_-$  and mirrored specific task de-  
357 mands. (3) Noise/exploration,  $\alpha_+$ , and  $\alpha_-$  all captured enough task-specific  
358 variance to be correctly identified with the corresponding task, showing that  
359 they captured different processes depending on the task (PC2-PC3).

### 360 *2.2.2. Parameters and Cognitive Processes*

361 Whereas the previous analysis revealed that all parameters contained  
362 task-specific information, it did not specify how much information was task-  
363 specific and how much was shared. For example, noise/exploration param-  
364 eters contained enough task-specific information to make it possible to de-  
365 termine in which task they were measured (PC2-PC3; Fig. 3B, 3C), but  
366 they also showed similar associations with engagement across tasks (PC1;  
367 Fig. 3A), similar age trajectories (Fig. 2C), and were mutually predictive  
368 (Fig. 2D). To quantify these patterns, we need to understand how much of  
369 each parameter's variance was unique and how much was shared between  
370 parameters and between tasks.

371 To achieve this, we probed how much of each parameter's variance was  
372 explained by other parameters, using regression. We assumed that param-  
373 eters reflected one or more cognitive processes, such that shared variance  
374 between parameters would imply overlapping cognitive processes. If param-  
375 eters reflected similar cognitive processes across tasks, then the same param-  
376 eter should dominate this analysis (e.g., when using parameters in task A  
377 to predict  $\frac{1}{\beta}$  in task B, task A's  $\frac{1}{\beta}$  should show the largest regression co-  
378 efficient). However, if parameters captured different processes across tasks,  
379 this would not be the case (e.g., all parameters of task A might predict  
380 task B's  $\frac{1}{\beta}$  equally). We used repeated, k-fold cross-validated Ridge regres-  
381 sion to avoid overfitting, obtaining unbiased out-of-sample estimates of the

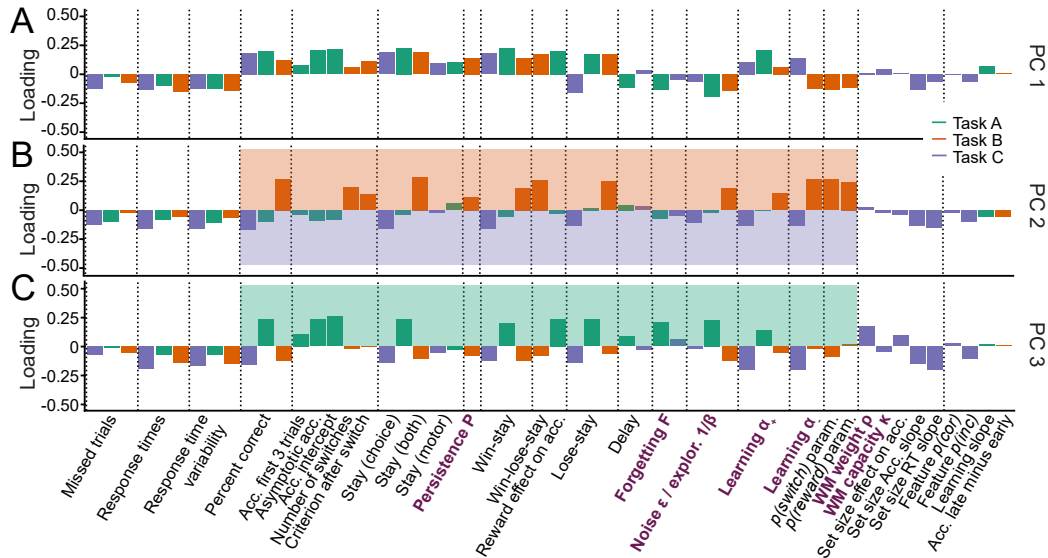


Figure 3: **Identifying the major axes of variation in the dataset.** A PCA was conducted on the entire dataset (39 behavioral features and 15 model parameters). The figure shows the factor loadings of the first three PCs. RL model parameters are highlighted in purple on the x-axis. Behavioral features are explained in detail in Appendix A and Appendix B. (A) PC1 captured broadly-defined task engagement, with negative loadings on features that were negatively associated with performance (e.g., number of missed trials) and positive loadings on features that were positively associated with performance (e.g., percent correct trials). (B-C) PC2 (B) and PC3 (C) captured task contrasts. PC2 loaded positively on features of task B (orange box) and negatively on features of task C (blue box). PC3 loaded positively on features of task A (green box) and negatively on features of tasks B and C. We flipped the loadings of features that were negative on PC1 when showing PC2 and PC3 to better visualize the task contrasts (section 4.6).

382 means and variances of explained variance  $R^2$  and regression coefficients  $w$   
383 (for methods, see section 4.7).

384 We first assessed the overall patterns of prediction, and found that all  
385 significant coefficients highlighted shared variance between tasks A and B or  
386 tasks A and C, but never between tasks B and C, mirroring our previous  
387 results (Fig. 2D; section 2.1.2) and patterns of task similarity (Fig. 1B).  
388 This means that no parameters in tasks B or C played a significant role in  
389 predicting parameters in the other, while both tasks' parameters were predic-  
390 tive (and being predicted by) parameters in task A. This further highlights  
391 the potential role of task similarity in parameter generalizability.

392 We next focused on noise/exploration parameters. Noise/exploration  
393 parameters in tasks B and C showed significant coefficients when predict-  
394 ing noise/exploration in task A, but the inverse was not true, such that  
395 noise/exploration in task A did not show significant coefficients when pre-  
396 dicting noise/exploration in tasks B or C (Fig. 4A; Table 5). The first  
397 result shows that noise/exploration parameters captured variance (cognitive  
398 processes) in task A that they also captured in tasks B and C. The second re-  
399 sult shows that noise/exploration parameters captured additional cognitive  
400 processes in tasks B and C that they did not capture in task A. Further-  
401 more, prediction accuracy increased when combining tasks B and C's param-  
402 eters to predict noise/exploration in task A, showing that noise/exploration  
403 parameters in tasks B and C captured partly non-overlapping aspects of  
404 noise/exploration in task A (Fig. 4B, left-most set of bars, compare pur-  
405 ple to orange and blue). This highlights both specificity in terms of which  
406 cognitive processes were captured by noise/exploration parameters across  
407 tasks (prediction between similar tasks), and some lack thereof (prediction  
408 was just one-way; no prediction between dissimilar tasks). Furthermore,  
409 noise/exploration in task A was predicted by Persistence and  $\alpha_-$  in task B,  
410 and by  $\alpha_-$  and working-memory weight  $\rho$  in task C (Fig. 4A; Table 5). This  
411 shows that some processes that noise/exploration parameters captured in  
412 task A were captured by different parameters in the other tasks, revealing a  
413 lack of distinctiveness in noise/exploration parameters.

414 We next assessed learning rates. Specificity was evident in that learning  
415 rate  $\alpha_+$  in task A showed a significant regression coefficient when predicting  
416 learning rates  $\alpha_+$  and  $\alpha_-$  in task C, and learning rate  $\alpha_-$  in task C showed  
417 a significant coefficient when predicting learning rate  $\alpha_+$  in task A (Fig. 4A;  
418 Table 5). However, a lack of specificity was evident in task B: When predict-  
419 ing  $\alpha_+$  in task B, no parameter of any task showed a significant coefficient

420 (including  $\alpha_+$  in other tasks; Table 5), and it was impossible to predict vari-  
421 ance in task B's  $\alpha_+$  even when combining all parameters of the other tasks  
422 (Fig. 4B, "Task B" panel). This reveals that  $\alpha_+$  captured fundamentally  
423 different cognitive processes in task B compared to the other tasks. The case  
424 was similar for parameter  $\alpha_-$ , which strikingly was inversely related between  
425 tasks A and B (Table 5), and impossible to predict in task B from all other  
426 parameters (Fig. 4B). This shows a lack of specificity, implying that learning  
427 rates did not reflect a consistent core of cognitive processes across tasks.

428 We then turned to the distinctiveness of learning rate parameters. Learn-  
429 ing rate  $\alpha_+$  in task A was predicted indistinctly by all parameters of task B  
430 (with the notable exception of  $\alpha_+$  itself; Fig. 4A; Table 5), suggesting that  
431 the cognitive processes that  $\alpha_+$  captured in task A were captured by an inter-  
432 play of several parameters in task B. Furthermore, task A's  $\alpha_+$  was predicted  
433 by task C's working-memory parameters  $\rho$  and  $K$  (Fig. 4A; Table 5), suggest-  
434 ing that  $\alpha_+$  captured a conglomerate of RL and working-memory processes  
435 in task A that was isolated by different sets of parameters in task C [52].  
436 In support of this interpretation, no variance in task C's working-memory  
437 parameters could be explained by any other parameters (Fig. 4B), revealing  
438 that they captured unique cognitive processes, likely working memory. Task  
439 C's RL parameters, on the other hand, could be explained by parameters in  
440 other tasks (Fig. 4B), suggesting they captured overlapping RL processes.

### 441 2.2.3. Parameters and Behavior

442 Faced with mounting evidence for parameter inconsistencies, we lastly  
443 aimed to uncover whether parameters shared any consistent similarities across  
444 tasks. The previous sections showed that parameters likely captured differ-  
445 ent (neuro)cognitive processes across tasks (e.g., different internal character-  
446 istics of learning and choice). However, computational models are funda-  
447 mentally models of behavior, so we argued that parameters might capture  
448 similar behavioral features (e.g., similar tendencies to stay after positive feed-  
449 back). Even though related, (neuro)cognitive processes and behavioral pat-  
450 terns should not be equated (Fig. 5). For example, different (neuro)cognitive  
451 mechanisms (e.g., prefrontal cortical reasoning, basal ganglia value learning,  
452 hippocampal episodic memory) might underlie the same behavioral pattern  
453 (e.g., lose-stay behavior) in different tasks, depending on the characteris-  
454 tics (e.g., stable versus volatile contingencies; deterministic versus stochastic  
455 feedback).

456 To investigate this possibility, we assessed the relationships between model

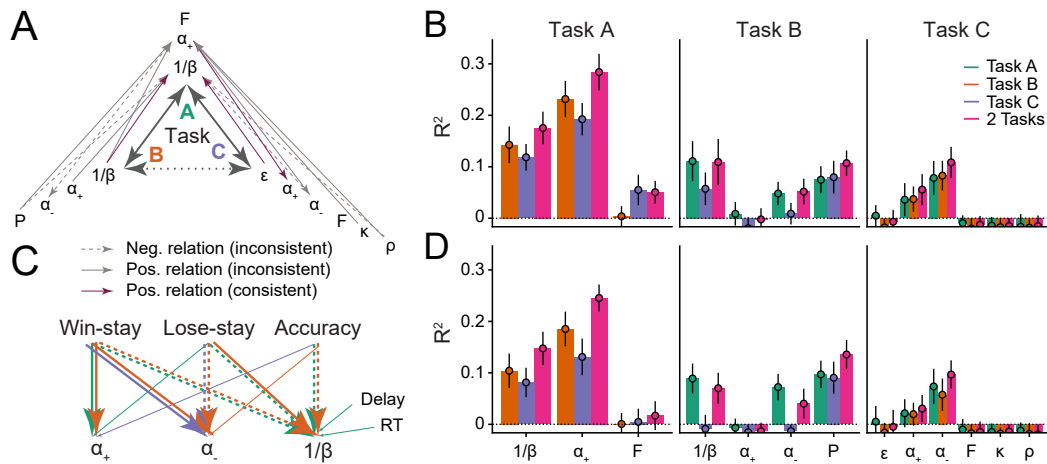


Figure 4: **Assessing parameter interpretability by analyzing shared variance.** (A) Parameter variance (cognitive processing) that is shared between tasks. Each arrow shows a significant regression coefficient when predicting a parameter in one task (e.g.,  $\alpha_+$  in task A) from all parameters of a different tasks (e.g.,  $P$ ,  $\alpha_-$ ,  $\alpha_+$ , and  $\frac{1}{\beta}$  in task B). The predicted parameter is shown at the arrow head, predictors at its end. Full lines indicate positive regression coefficients, and are highlighted in purple when connecting two identical parameters; dotted lines indicate negative coefficients; non-significant coefficients are not shown. Table 5 provides the full statistics of the models summarized in this figure. (B) Variance of each parameter that was also captured by parameters of other models. Each bar shows the percentage of explained variance ( $R^2$ ) when predicting one parameter from all parameters of a different task/model, using Ridge regression. Part (A) of this figure shows the coefficients of these models. The x-axis shows the predicted parameter, and colors differentiate between predicting tasks. Three models were conducted to predict each parameter: One combined the parameters of both other tasks (pink), and two kept them separate (green, orange, blue). Larger amounts of explained variance (e.g., Task A  $\frac{1}{\beta}$  and  $\alpha_-$ ) suggest more shared processes between predicted and predicting parameters; the inability to predict variance (e.g., Task B  $\alpha_+$ ; Task C working-memory parameters) suggests that distinct processes were captured. Bars show mean  $R^2$ , averaged over  $k$  data folds ( $k$  was chosen for each model based on model fit, using repeated cross-validated Ridge regression; for details, see section 4.7); error bars show standard errors of the mean across folds. (C) Relations between parameters and behavior. The arrows visualize Ridge regression models that predict parameters (bottom row) from behavioral features (top row) within tasks (full statistics in Table 6). Arrows indicate significant regression coefficients, colors denote tasks, and line types denote the sign of the coefficients, like before. All significant within-task coefficients are shown. Task-based consistency (similar relations between behaviors and parameters across tasks) occurs when arrows point from the same behavioral features to the same parameters in different tasks (i.e., multiple arrows). (D) Variance of each parameter that was explained by behavioral features; corresponds to the behavioral Ridge models shown in part (C).

457 parameters and behavioral features across tasks. Using regularized Ridge re-  
458 gression like above, we predicted each model parameter from five selected  
459 behavioral features (Appendix A, Appendix C) of each of the three tasks  
460 (15 predictors; for regression methods, see section 4.7). One possible outcome  
461 of this analysis is “absolute consistency”: parameters might capture the same  
462 behavioral pattern within and across tasks (e.g., noise/exploration of each  
463 task might capture task A accuracy). This outcome would be expected if pa-  
464 rameters captured the same cognitive processes across tasks, and behavioral  
465 features were a direct reflection of cognitive processes. Another possible out-  
466 come is “absolute *in*consistency” (e.g., in every task, noise/exploration might  
467 capture different behavioral features). This outcome would suggest that pa-  
468 rameters captured unrelated cognitive and behavioral features in each task.  
469 Crucially, a third possible outcome is “task-based consistency”: Parameters  
470 might capture the same behavioral features, but only within tasks (e.g., in  
471 each task, learning rates might capture the win-stay behavior of that task,  
472 but not of other tasks). This outcome would suggest that parameters gen-  
473 eralized in terms of which behavioral features they reflected, but behavioral  
474 features—like (neuro)cognitive processes—differed between tasks.

475 Focusing on noise/exploration parameters,  $\frac{1}{\beta}$  in tasks A and B was pre-  
476 dicted by task A win-stay behavior, revealing absolute consistency (Table  
477 6).  $\frac{1}{\beta}$  was also predicted by accuracy, win-stay, and lose-stay behavior within  
478 both tasks A and B, but not across tasks, revealing task-based consistency  
479 (Fig. 4C; Table 6). For learning rates,  $\alpha_+$  in tasks A and B was predicted  
480 by the corresponding win-stay behavior, and  $\alpha_-$  in tasks B and C was neg-  
481 atively predicted by the corresponding lose-stay behavior, and positively by  
482 the corresponding win-stay behavior (Fig. 4C; Table 6), revealing task-based  
483 consistency. The consistency of  $\alpha_-$  is especially noteworthy given the abun-  
484 dance of discrepancies in previous sections.

485 Taken together, noise/exploration parameters,  $\alpha_+$ , and  $\alpha_-$  captured simi-  
486 lar behavioral features across tasks (Fig. 4C), despite differences in cognitive  
487 processing (Fig. 4A, 4B), captured information (Fig. 3B, 3C), age trajec-  
488 tories (Fig. 2C, 2D), and absolute values (Fig. 2A, 2B). Notably, the observed  
489 discrepancies reflected task characteristics (Fig. 1B, 3A) for both param-  
490 eters (Fig. 1B, 3A) and behavior (suppl. Fig. B.6B), suggesting that task  
491 characteristics shaped behavioral responses and model parameters.

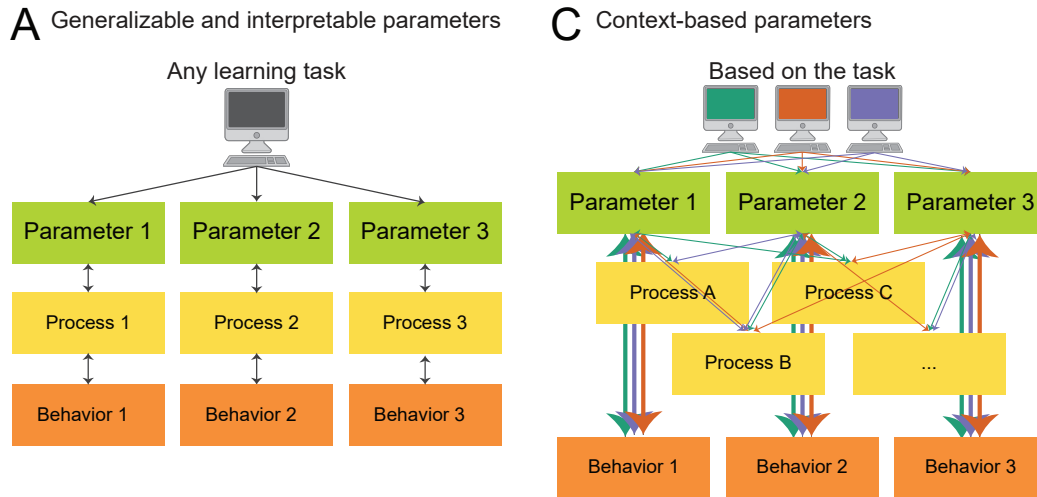


Figure 5: **What do model parameters measure?** (A) View based on generalizability and interpretability. In this view, which is implicitly taken by much current computational modeling research, models are fitted in order to reveal individuals' intrinsic model parameters, which reflect clearly delineated, separable, and meaningful (neuro)cognitive processes, a concept we call *interpretability*. Interpretability is evident in that every model parameter captures a specific cognitive process (bidirectional arrows between each parameter and process), and that cognitive processes are separable from each other (no connections between processes). Task characteristics are treated as irrelevant, a concept we call *generalizability*, such that parameters of any learning task (within reason) are expected to capture similar cognitive processes. (C) Updated view, based on our results, that acknowledges the role of context (e.g., task characteristics, model parameterization, participant sample) in computational modeling. Which cognitive processes are captured by each model parameter is influenced by context (green, orange, blue), as shown by distinct connections between parameters and cognitive processes. Different parameters within the same task can capture overlapping cognitive processes (not interpretable), and the same parameters can capture different processes depending on the task (not generalizable). However, parameters likely capture consistent behavioral features across tasks (thick vertical arrows).



### 492 3. Discussion

493 The current study subjected the generalizability and interpretability of  
494 RL models to a scrupulous empirical investigation, using a developmental  
495 sample. We found weaker levels of generalizability and interpretability than  
496 would be expected based on current research practices, such as comparing  
497 parameters between studies that use different tasks and models [28].

498 Interestingly, patterns of generalizability and interpretability varied be-  
499 tween parameters: Exploration/noise parameters showed considerable gener-  
500 alizability in the form of correlated variance and age trajectories. The decline  
501 in exploration/noise we observed between ages 8-17 was also consistent with  
502 previous studies reviewed in [13]. Interpretability of exploration/noise pa-  
503 rameters was mixed: Despite evidence for specificity in some cases (overlap  
504 in parameter variance between some tasks), it was missing in others (lack of  
505 overlap between other tasks), and crucially, parameters lacked distinctiveness  
506 (substantial overlap in variance with other parameters).

507 Learning rate from negative feedback, on the other hand, showed a sub-  
508 stantial lack of generalizability: parameters were less consistent within par-  
509 ticipants than within tasks, and age trajectories differed both quantitatively  
510 and qualitatively. This result is consistent with discrepancies in learning rate  
511 parameters across developmental studies [13]. Learning rates from positive  
512 and negative feedback combined were interpretable to a limited degree (over-  
513 lap in variance between some tasks). However, interpretability was overshad-  
514 owed by a lack of specificity (lack of shared core variance) and distinctive-  
515 ness (fundamental entangling with several other parameters, most notably  
516 working-memory parameters).

517 These within-participant findings are consistent with patterns that are  
518 emerging from comparisons of studies published by different labs [13]. Our  
519 within-participant design allowed us to go beyond these between-study find-  
520 ings by confirming that the same participants can show different parame-  
521 ters when tested using different tasks. The within-participant consistency of  
522 noise/exploration parameters strengthens our confidence that these indeed  
523 decrease with age [13, 53, 54]. The inconsistency of learning rate parameters  
524 leads to the unexpected, but important conclusion that we cannot measure  
525 an individual’s “*intrinsic* learning rate” using RL modeling, and that we  
526 cannot draw general conclusions about “the development of learning rates”  
527 that apply to all RL contexts, using current methods.

528 Our findings also help us clarify the source of parameter inconsistencies

529 in the previous literature, which could indicate replication problems and  
530 technical issues: For example, model misspecification [13], lack of model  
531 comparison and validation [27, 30], inappropriate fitting methods [29, 51],  
532 and lack of parameter reliability due to suboptimal methods [48] have all  
533 been suggested as potential sources of inconsistencies. However, our results  
534 show that discrepancies are expected even with a consistent methodological  
535 pipeline, and using up-to-date modeling techniques (detailed model compari-  
536 son, validation, and hierarchical Bayesian model fitting where possible). This  
537 should encourage the field of computational modeling to study the external  
538 factors that drive such inconsistencies, and are currently undescribed by RL  
539 methods, with more rigor.

### 540 *3.1. Limitations*

541 One limitation of our results is that regression analyses might be con-  
542 taminated by parameter cross-correlations (in sections 2.1.2, 2.1.3, 2.1.4),  
543 which would reflect modeling limitations (fewer true degrees of freedom than  
544 model parameters), and not necessarily shared cognitive processes. For ex-  
545 ample, parameters  $\alpha$  and  $\beta$  are mathematically related in the regular RL  
546 modeling framework [26, 29], and we observed significant correlations be-  
547 tween parameters within tasks for two of our three tasks (suppl. Fig. D.7).  
548 This indicates that caution is required when interpreting correlation results.  
549 However, correlations were also present between tasks (suppl. Fig. D.8),  
550 suggesting that within-model trade-offs were not the only explanation for  
551 shared variance, and that shared cognitive processes likely also played a role.  
552 Furthermore, correlations between parameters within models are frequent  
553 in the existing literature, and do not prevent researchers from interpreting  
554 parameters—in this sense, the existence of similar correlations in our study  
555 allows us to address the question of generalizability and interpretability in  
556 similar circumstances as in the existing literature.

### 557 *3.2. Moving Forward*

558 With this research, we do not intend to undermine RL modeling, but to  
559 improve its quality. Computational model parameters potentially provide  
560 highly valuable insights into (neuro)cognitive processing—we just need to  
561 refrain from assuming that the identified processes are necessarily and inher-  
562 ently specific, distinct, and “theoretically meaningful” [14] (interpretable).  
563 Parameters with the same names also do not automatically transfer between

564 tasks or models, and are less interchangeable than we often implicitly as-  
565 sume [28]. At the same time, the behavioral features that are captured by  
566 parameters seem to generalize well between tasks.

567 In the long term, we need to understand why RL parameters differ be-  
568 tween tasks. We suggest three potential, not mutually exclusive answers:

569 1. **Optimality.** Variance in RL parameters may reflect how participants  
570 adapt their behavior to task demands, an explanation proposed by  
571 [13]. For example, participants might tune learning rates to task char-  
572 acteristics (e.g., adopting lower learning rates in stable than volatile  
573 contexts [55]), rather than learning rates reflecting intrinsic “settings”  
574 (e.g., 10-year-olds having a learning rate of 20%; 16-year-olds of 40%).  
575 An optimality-based view would also explain why learning rates dif-  
576 fer between deterministic and stochastic tasks, which require different  
577 amounts of behavioral change in response to feedback, to reach opti-  
578 mal performance. Age differences can potentially be explained because  
579 optimal settings likely differ between ages because they interact with  
580 different environments, or because different ages might have different  
581 capacities to shift internal settings when shifting from task to task.  
582 More research is needed, however, to determine whether parameter opti-  
583 mality can explain all inconsistencies in the literature. For example,  
584 our finding that participants showed the most optimal parameter val-  
585 ues in the intermediate age range in task B [49], whereas optimality  
586 increased monotonously with age in tasks A and C [35, 50], is difficult  
587 to reconcile with this view.

588 2. **Modulatory processes.** RL Parameters may vary as a function of  
589 modulatory processes that are not well captured in current RL models.  
590 Modulatory processes have been described in cognition and neurobiol-  
591 ogy and likely serve to shift functional outputs (e.g., hunger increasing  
592 motivation) [56, 57, 58]. Some modulatory processes reflect the inte-  
593 gration of external contextual information: for example, uncertainty  
594 affects dopamine neuron firing [59, 60, 61]. In addition, environments  
595 with different degrees of uncertainty have been shown to elicit different  
596 learning rates [55]. It is thus possible that neuromodulation by task  
597 uncertainty could modulate RL processes, reflected in RL parameters.  
598 In our data, feedback stochasticity and task volatility likely contribute  
599 to such uncertainty-related modulation. However, other factors like  
600 task similarity (low versus high), task characteristics (e.g., volatility

601 [55, 49], feedback stochasticity, memory load [35, 52], feedback valence  
602 and conditioning type [24]), and choice of model parameters (e.g., for-  
603 getting [35, 50], counter-factual learning [49], negative and positive  
604 learning rates [34, 62, 63]), also seem to affect RL parameters, but are  
605 independent of uncertainty. More research is needed to systematically  
606 investigate the factors that contribute to modulatory processes, and how  
607 they impact cognition and computation.

608 **3. RL processes are multifaceted.** RL Parameters capture a multi-  
609 tude of separate processes, whose composition differs across tasks (Fig.  
610 5B). RL algorithms are framed in the most general way to allow ap-  
611 plication to a wide range of contexts, including AI, neuroscience, and  
612 psychology [26, 28]. As behavioral models, their use has spanned be-  
613 haviors from simple conditioning [1, 38] to complex decision making  
614 [4, 6, 7, 8, 9, 22, 64], meaning that the same parameters capture cog-  
615 nitive processes that vary considerably in type and complexity: Pro-  
616 cesses can include the slow acquisition of implicit preferences [1], long-  
617 term memory for such preferences [65], quick recognition of contingency  
618 switches [49, 66], selection of abstract high-level strategies [7, 9, 67],  
619 meta-learning [6], habitual and goal-directed decision making [5], work-  
620 ing memory or episodic memory-guided choice [52, 68, 69], and many  
621 others. This list alone outnumbers the list of typical RL model param-  
622 eters, suggesting that RL parameters capture different (combinations  
623 of) cognitive processes depending on the paradigm. Similar arguments  
624 have also been made for behavioral analyses [70].

### 625 *3.3. Conclusion*

626 Our research has important implications for fields that focus on individ-  
627 ual differences, including developmental and clinical computational research.  
628 The current study should be seen as a proof of concept that many contextual  
629 factors impact computational modeling, and larger studies will be necessary  
630 to quantify these effects and determine their structure. Other areas of model-  
631 ing besides the RL framework should be subjected to a similar investigation.  
632 It is possible, for example, that generalizability differs for sequential sam-  
633 pling [71, 72], Bayesian inference [49, 73, 74], model-based versus model-free  
634 RL [48, 75, 76], or other models.

635 In sum, our results suggest that relating model parameters to cognitive  
636 constructs and real-world behavior might require us to carefully account for  
637 task variables, and environmental variability in general. This ties into the

638 bigger picture of understanding how neurocognitive processes are shared be-  
639 tween tasks [77], and reflects a larger pattern of realization in psychology  
640 that we cannot objectively assess an individual’s cognitive processing while  
641 ignoring subjective context. We have shown that in lab studies, different task  
642 contexts recruit different system settings within an individual; similarly, real-  
643 life environment, its changes during development, and past environment [78]  
644 may also modulate which cognitive processes we recruit. Heightened aware-  
645 ness and systematic study of contextual variables will therefore be a valuable  
646 future investment as we work to measure and accommodate diversity in cog-  
647 nitive processes.

## 648 4. Methods

### 649 4.1. Study Design

650 Our sample of 291 participants was balanced between females and males,  
651 and all ages (8-30 years) were represented equally (Fig. 1A, left). Partici-  
652 pants completed four computerized tasks, questionnaires, and a saliva sample  
653 during the 1-2 hour lab visit (see section 4.3). To reduce noise, we excluded  
654 participants based on task-specific performance criteria (see section 4.2). Due  
655 to worse performance, more younger than older participants were excluded,  
656 which is a caveat for the interpretation of age effects (note however that  
657 these exclusions cannot account for the observed age effects but act against  
658 them; Fig. 1A). Our tasks—A (“Butterfly task” [50, 79]), B (“Probabilis-  
659 tic Switching” [66, 49]), and C (“Reinforcement learning-Working memory”  
660 [35, 52])—were all classic reinforcement learning tasks: on each trial, partic-  
661 ipants chose between several actions in an effort to earn rewards, which were  
662 presented as binary feedback (win/point or lose/no point) after each choice.

663 The tasks varied on several common dimensions (Fig. 1B), which have  
664 been related to discrepancies in behavioral and neurocognitive results in the  
665 literature [24, 36, 37]. For example, in one task (task C), positive feedback  
666 was deterministic, such that every correct action led to a positive outcome,  
667 whereas in the two other tasks (tasks A and B), positive feedback was stochas-  
668 tic, such that some correct actions led to positive and others to negative  
669 outcomes. A different set of two tasks (B and C) provided diagnostic posi-  
670 tive feedback, such that every positive outcome indicated a correct action,  
671 whereas in the third (A), positive feedback was non-diagnostic, such that  
672 positive outcomes could indicate both correct and incorrect actions. Two  
673 tasks (A and C) presented several different stimuli/states for which correct

674 actions had to be learned, whereas the third (B) only presented a single  
675 one. Overall, task A shared more important similarities with both tasks B  
676 and C than either of these shared with each other, allowing us to explore  
677 whether task similarity played a role in parameter generalizability and in-  
678 terpretability. A comprehensive list of task differences is shown in Fig. 1B,  
679 and each task is described in more detail in section 4.4. Section Appendix  
680 B explains the most prominent findings of each task individually, and shows  
681 several behavioral measures over age.

## 682 *4.2. Participant Sample*

### 683 *4.2.1. Sample Overview*

684 All procedures were approved by the Committee for the Protection of Hu-  
685 man Subjects at the University of California, Berkeley. We tested 312 partic-  
686 ipants: 191 children and adolescents (ages 8-17) and 55 adults (ages 25-30)  
687 were recruited from the community and completed a battery of computer-  
688 ized tasks, questionnaires, and saliva samples; 66 university undergraduate  
689 students (aged 18-50) completed the four tasks as well, but not the question-  
690 naires or saliva sample. Community participants of all ages were prescreened  
691 for the absence of present or past psychological and neurological disorders;  
692 the undergraduate sample indicated the absence of these. Compensation for  
693 community participants consisted in \$25 for the 1-2 hour in-lab portion of  
694 the experiment and \$25 for completing optional take-home saliva samples;  
695 undergraduate students received course credit for participation in the 1-hour  
696 study.

### 697 *4.2.2. Participant Exclusion*

698 Two participants from the undergraduate sample were excluded because  
699 they were older than 30, and 7 were excluded because they failed to indicate  
700 their age. This led to a sample of 191 community participants under 18, 57  
701 undergraduate participants between the ages of 18-28, and 55 community  
702 participants between the ages of 25-30. Of the 191 participants under 18,  
703 184 completed task B, and 187 completed tasks A and C. Reasons for not  
704 completing a task included getting tired, running out of time, and technical  
705 issues. All 57 undergraduate participants completed tasks B and C and  
706 55 completed task A. All 55 community adults completed tasks B and A,  
707 and 45 completed task C. Appropriate exclusion criteria were implemented  
708 separately for each task to exclude participants who failed to pay attention  
709 and who performed critically worse than the remaining sample (for task A,



710 see [50]; task B [49]; task C [35]). Based on these criteria, 5 participants  
711 under the age of 18 were excluded from task B, 10 from task A, and none  
712 from task C. One community adult participant was excluded from task A,  
713 but no adult undergraduates or community participants were excluded from  
714 tasks B or C.

715 Because this study related the results across all three tasks, we only  
716 included participants who were not excluded in any task, leading to a final  
717 sample of 143 participants under the age of 18 (male: 77; female: 66), 51  
718 undergraduate participants (male: 17; female: 34), and 53 adults from the  
719 community (male: 25; female: 28), for a total of 247 participants (male: 119;  
720 female: 128). We entirely excluded the fourth task of our study from the  
721 current analysis, which was modeled after a rodent task and used in humans  
722 for the first time [80], because the applied performance criterion led to the  
723 exclusion of the majority of our developmental sample. We split participants  
724 into quantiles based on age, which were calculated separately within each  
725 sex.

#### 726 *4.3. Testing Procedure*

727 After entering the testing room, participants under 18 years and their  
728 guardians provided informed assent and permission; participants over 18  
729 provided informed consent. Guardians and participants over 18 filled out  
730 a demographic form. Participants were led into a quiet testing room in view  
731 of their guardians, where they used a video game controller to complete  
732 four computerized tasks. The first task was called “4-choice” and assessed  
733 reversal learning in an environment with 4 different choice options, with a  
734 duration of approximately 5 minutes (designed after [80]). This task was  
735 excluded from the current analysis (see section 4.2.2). The second task was  
736 C (“Reinforcement learning-Working memory”) and took about 25 minutes  
737 to complete [52, 35]. After the second task, participants between the ages of  
738 8-17 provided a saliva sample (for details, see [35]) and took a snack break (5-  
739 10 minutes). After that, participants completed task A (“Butterfly task”),  
740 which took about 15 minutes [79, 50], and task B (“Probabilistic Switch-  
741 ing”), which took about 10 minutes to complete [49]. At the conclusion of  
742 the tasks, participants between 11 and 18 completed the Pubertal Develop-  
743 ment Scale (PDS [81]) and were measured in height and weight. Participants  
744 were then compensated with \$25 Amazon gift cards. The PDS questionnaire  
745 and saliva samples were administered to investigate the role of pubertal mat-  
746 uration on learning and decision making. Pubertal analyses are not the focus



747 of the current study and will be or have reported elsewhere [35, 49, 50]. For  
748 methodological details, refer to [35]. The entire lab visit took 60-120 minutes,  
749 depending on the participant.

#### 750 *4.4. Task Design*

##### 751 *4.4.1. Task A (“Butterfly task”)*

752 The goal of task A was to collect as many points as possible, by guessing  
753 correctly which of two flowers was associated with each of four butterflies.  
754 Correct guesses were rewarded with 70% probability, and incorrect guesses  
755 with 30%. The task contained 120 trials (30 for each butterfly) that were  
756 split into 4 equal-sized blocks, and took between 10-20 minutes to complete.  
757 More detailed information about methods and results can be found in [50].

##### 758 *4.4.2. Task B (“Probabilistic Switching”)*

759 The goal of task B was to collect golden coins, which were hidden in two  
760 green boxes. The task could be in one of two states: “Left box is correct”  
761 or “Right box is correct”. In the former, selecting the left box led to reward  
762 in 75% of trials, while selecting the right box never led to a reward (0%).  
763 Several times throughout the task, task contingencies changed unpredictably  
764 and without notice (after participants had reached a performance criterion  
765 indicating they had learned the current state), and the task switched states.  
766 Participants completed 120 trials of this task (2-9 reversals), which took ap-  
767 proximately 5-15 minutes. For more information and additional task details,  
768 refer to [49].

##### 769 *4.4.3. Task C (“Reinforcement Learning-Working Memory”)*

770 The goal of task C was to collect as many points as possible by pressing  
771 the correct key for each stimulus. Pressing the correct key deterministically  
772 led to reward, and the correct key for a stimulus never changed. Stimuli  
773 appeared in blocks that varied in the number of different stimuli, with set  
774 sizes ranging from 2-5. In each block, each stimulus was presented 12-14  
775 times, for a total of 13 \* set size trials per block. Three blocks were presented  
776 for set sizes 2-3, and 2 blocks were presented for set sizes 4-5, for a total of 10  
777 blocks. The task took between 15-25 minutes to complete. For more details,  
778 as well as a full analysis of this dataset, refer to [35].

779 *4.5. Computational Models*

For all tasks, we used RL theory to model how participants adapted their behavior in order to maximize reward. RL models assume that agents learn a policy  $\pi(a|s)$  that determines (probabilistically) which action  $a$  to take in each state  $s$  of the world [26]. Here and in most cognitive RL models, this policy is based on action values  $Q(a|s)$ , i.e., the values of each action  $a$  in each state  $s$ . Agents learn action values by observing the reward outcomes,  $r_t$ , of their actions at each time step  $t$ . Learning consists in updating existing action values  $Q_t(a|s)$  using the “reward prediction error”, the difference between the expected reward  $Q_t(a|s)$  and the actual reward  $r_t$ :

$$Q_{t+1}(a|s) = Q_t(a|s) + \alpha(r_t - Q_t(a|s))$$

780 How much a learner weighs past action value estimates compared to new out-  
781 comes is determined by parameter  $\alpha$ , the learning rate. Small learning rates  
782 favor past experience and lead to stable learning over long time horizons,  
783 while large learning rates favor new outcomes and allow for faster and more  
784 flexible changes, focusing on shorter time horizons. With enough time and  
785 in a stable environment, the RL updating scheme guarantees that value es-  
786 timates will reflect the environment’s true reward probabilities, and thereby  
787 allow for optimal long-term choices [26].

In order to choose actions, most cognitive RL models use a (noisy) “soft-  
max” function to translate action values  $Q(a|s)$  into policies  $p(a|s)$ :

$$p(a_i|s) = \frac{\exp(\beta Q(a_i|s))}{\sum_{a_j \in A} \exp(\beta Q(a_j|s))}$$

788  $A$  refers to the set of all available actions (tasks A and B have 2 actions,  
789 task C has 3), and  $a_i$  and  $a_j$  to individual actions within the set. How  
790 deterministically versus noisily this translation is executed is determined by  
791 exploration parameters  $\beta$ , also called inverse decision temperature, and/or  
792  $\epsilon$ , the decision noise (see below). Small decision temperatures  $\frac{1}{\beta}$  favor the  
793 selection of the highest-valued actions, enabling exploitation, whereas large  
794 decision temperatures select actions of low and high values more evenly,  
795 enabling exploration. Parameter  $\epsilon$  adds undirected noise to action selection,  
796 selecting random action with a small probability  $\epsilon$  on each trial.

797 Besides  $\alpha$ ,  $\beta$ , and noise, cognitive RL models often include additional  
798 parameters to better fit empirical behavior in humans or animals. Com-  
799 mon choices include Forgetting—a consistent decay of action values back to

800 baseline—, and Persistence—the tendency to repeat the same action inde-  
 801 pendent of outcomes, a parameter also known as sticky choice or perseverance  
 802 [63]. In addition, cognitive models often differentiate learning from positive  
 803 versus negative rewards, splitting learning rate  $\alpha$  into two separate param-  
 804 eters  $\alpha_+$  and  $\alpha_-$ , which are applied to only positive and only negative out-  
 805 comes, respectively [34, 40, 82, 83, 84, 85, 86, 87, 88]. The next paragraphs  
 806 introduce these parameters in detail.

In task A, the best fitting model included a forgetting mechanism, which was implemented as a decay in Q-values applied to all action values of the three stimuli (butterflies) that were not shown on the current trial:

$$Q_{t+1}(a|s) = (1 - f) * Q_t(a|s) + f * 0.5.$$

807 The free parameter  $0 < f < 1$  reflects individuals’ tendencies to forget.

In task B, free parameter  $P$  captured choice persistence, which biased choices on the subsequent trial toward staying ( $P > 0$ ) or switching ( $P < 0$ ).  $P$  modifies action values  $Q(a|s)$  into  $Q'(a|s)$ , as follows:

$$Q'_t(a|s) = Q_t(a|s) + P \iff a_t = a_{t-1}$$

$$Q'_t(a|s) = Q_t(a|s) \iff a_t \neq a_{t-1}$$

In addition, the model of task B included counter-factual learning parameters  $\alpha_{C+}$  and  $\alpha_{C-}$ , which added counter-factual updates based on the inverse outcome and affected the non-chosen action. For example, after receiving a positive outcome ( $r = 1$ ) for choosing left ( $a$ ), counter-factual updating would lead to an “imaginary” negative outcome ( $\bar{r} = 0$ ) for choosing right ( $\bar{a}$ ).

$$Q_{t+1}(\bar{a}|s) = Q_t(\bar{a}|s) + \alpha_{C+}(\bar{r} - Q_t(\bar{a}|s)) \iff r = 1$$

$$Q_{t+1}(\bar{a}|s) = Q_t(\bar{a}|s) + \alpha_{C-}(\bar{r} - Q_t(\bar{a}|s)) \iff r = 0$$

808  $\bar{a}$  indicates the non-chosen action, and  $\bar{r}$  indicates the inverse of the received  
 809 outcome,  $\bar{r} = 1 - r$ . The best model fits were achieved with  $\alpha_{C+} = \alpha_+$  and  
 810  $\alpha_{C-} = \alpha_-$ , so counter-factual learning rates are not reported in this paper.

In tasks A and B, positive and negative learning rates are differentiated in the following way:

$$Q_{t+1}(a|s) = Q_t(a|s) + \alpha_+(r_t - Q_t(a|s)) \iff r_t = 1$$

$$Q_{t+1}(a|s) = Q_t(a|s) + \alpha_-(r_t - Q_t(a|s)) \iff r_t = 0$$

811 In the best model for task A, only  $\alpha_+$  was a free parameter, while  $\alpha_-$  was  
 812 fixed to 0. In task C,  $\alpha_-$  was a function of  $\alpha_+$ , such that  $\alpha_- = b * \alpha_+$ , where  $b$   
 813 is the neglect bias parameter that determines how much negative feedback is  
 814 neglected compared to positive feedback. Throughout the paper, we report  
 815  $\alpha_- = b * \alpha_+$  for task C.

In addition to an RL module, the model of task C included a working-memory module with perfect recall of recent outcomes, but subject to forgetting and capacity limitations. Perfect recall was modeled as an RL process with learning rate  $\alpha_{WM+} = 1$  that operated on working-memory weights  $W(a|s)$  rather than action values. On trials with positive outcomes ( $r = 1$ ), the model reduces to:

$$W_{t+1}(a|s) = r_t$$

On trials with negative outcomes ( $r = 0$ ), multiplying  $\alpha_{WM+} = 1$  with the neglect bias  $b$  leads to potentially less-than perfect memory:

$$W_{t+1}(a|s) = W_t(a|s) + b * (r_t - W_t(a|s))$$

Working-memory weights  $W(a|s)$  were transformed into action policies  $p_{WM}(a|s)$  in a similar way as RL weights  $Q(a|s)$  were transformed into action probabilities  $p_{RL}(a|s)$ , using a softmax transform combined with undirected noise:

$$p(a_i|s) = (1 - \epsilon) * \frac{\exp(\beta Q(a_i|s))}{\sum_{a_j \in a} \exp(\beta Q(a_j|s))} + \epsilon * \frac{1}{|a|}$$

816  $|a| = 3$  is the number of available actions and  $\frac{1}{|a|}$  is the uniform policy over  
 817 these actions;  $\epsilon$  is the undirected noise parameter.

Forgetting was implemented as a decay in working-memory weights  $W(a|s)$  (but not RL Q-values):

$$W_{t+1}(a|s)_{t+1} = (1 - f) * W_t(a|s)_t + f * \frac{1}{3}$$

Capacity limitations of working memory were modeled as an adjustment in the weight  $w$  of  $p_{WM}(a|s)$  compared to  $p_{RL}(a|s)$  in the final calculation of action probabilities  $p(a|s)$ :

$$w = \rho * (\min(1, \frac{K}{ns}))$$

$$p(a|s) = w * p_{WM}(a|s) + (1 - w) * p_{RL}(a|s)$$

818 The free parameter  $\rho$  is the individual weight of working memory compared  
819 to RL,  $ns$  indicates a block’s stimulus set size, and  $K$  captures individual  
820 differences in working-memory capacity.

821 We fitted a separate RL model to each task, using state-of-the-art meth-  
822 ods for model construction, fitting, and validation [30, 27]. Models for tasks  
823 A and B were fitted using hierarchical Bayesian methods with Markov-Chain  
824 Monte-Carlos sampling, which is an improved method compared to maximum  
825 likelihood that leads to better parameter recovery, amongst other advantages  
826 [89, 90, 91]. The model for task C was fitted using classic non-hierarchical  
827 maximum-likelihood because model parameter  $K$  is discrete, which renders  
828 hierarchical sampling less tractable. In all cases, we verified that the model  
829 parameters were recoverable by the selected model-fitting procedure, and  
830 that the models were identifiable. Details of model-fitting procedures are  
831 provided in the original publications [35, 49, 50].

832 For additional details on any of these models, as well as detailed model  
833 comparison and validation, the reader is referred to the original publications  
834 (task A: [50]; task B: [49]; task C: [35]).

#### 835 4.6. Principal Component Analysis (PCA)

836 The PCA in section 2.2.1 included 15 model parameters ( $\alpha_+$  and noise/exploration  
837 in each task; Forgetting and  $\alpha_-$  in two tasks; Persistence in task B; four  
838 working-memory parameters in task C; see section 4.5) and 39 model-free  
839 features, including simple behavioral features (e.g., overall performance, re-  
840 action times, tendency to switch), results of behavioral regression models  
841 (e.g., effect of stimulus delay on accuracy), and the model parameters of an  
842 alternative Bayesian inference model in task B. All behavioral features, in-  
843 cluding their development over age, are described in detail in Appendix C  
844 and suppl. Fig. B.6B. For simplicity, section 2.2.1 focused on the first three  
845 PCs only; the weights, explained variance, and age trajectories of remaining  
846 PCs are shown in suppl. Fig. D.10.

847 PCA is a statistical tool that decomposes the variance of a dataset into  
848 so-called “principal components” (PCs). PCs are linear combinations of a  
849 dataset’s original features (e.g., response times, accuracy, learning rates),  
850 and explain the same variance in the dataset as the original features. The  
851 advantage of PCs is that they are orthogonal to each other and therefore  
852 capture independent aspects of the data. In addition, subsequent PCs ex-  
853 plain subsequently less variance, such that selecting just the top PCs of a  
854 dataset retains the bulk of the variance and the ability to reconstruct the

855 dataset up to random noise. When using this approach, it is important to  
856 understand which concept each PC captures. So-called factor loadings, the  
857 original features' weights on each PC, can provide this information.

858 PCA performs a *change of basis*: Instead of describing the dataset using  
859 the original features (in our case, 54 behaviors and model parameters), it cre-  
860 ates new features, PCs, that are linear combinations of the original features  
861 and capture the same variance, but are orthogonal to each other. PCs are  
862 created by eigendecomposition of the covariance matrix of the dataset: the  
863 eigenvector with the largest eigenvalue shows the direction in the dataset in  
864 which most variance occurs, and represents the first PC. Eigenvectors with  
865 subsequently smaller eigenvalues form subsequent PCs. PCA is related to  
866 Factor analysis, and often used for dimensionality reduction. In this case,  
867 only a small number of PCs is retained whereas the majority is discarded, in  
868 an effort to retain most variance with a reduced number of features.

869 We highlight the most central behavioral features here; more detail is pro-  
870 vided in Appendix A and Appendix C. Response to feedback was assessed  
871 using features “Win-stay” (percentage of trials in which a rewarded choice  
872 was repeated), and “Lose-stay” (percentage of trials in which a non-rewarded  
873 choice was repeated). For task B, we additionally included “Win-lose-stay”  
874 tendencies, which is the proportion of trials in which participants stay after a  
875 winning trial that is followed by a losing trial. This is an important measure  
876 for this task because the optimal strategy required staying after single losses.

877 We also included behavioral persistence measures in all tasks. In tasks  
878 A and C, these included a measure of action repetition (percentage of trials  
879 in which the previous key was pressed again, irrespective of the stimulus  
880 and feedback) and choice repetition (percentage of trials in which the action  
881 was repeated that was previously selected for the same stimulus, irrespective  
882 of feedback). In task B, both measures were identical because every trial  
883 presents the same stimulus.

884 We further included task-specific measures of performance. In task A,  
885 these were: the average accuracy for the first three presentations of each  
886 stimulus, reflecting early learning speed; and the asymptote, intercept, and  
887 slope of the learning progress in a regression model predicting performance  
888 (for details about these measures, see [50]). In task B, task-specific mea-  
889 sures of performance included the number of reversals (because reversals  
890 were performance-based); and the average number of trials to reach criterion  
891 after a switch. In tasks A and C, we also included a model-independent  
892 measure of forgetting. In task A, this was the effect of delay on performance

893 in the regression model mentioned above. In task C, this was the effect of  
894 delay in a similar regression model, which also included set size, the number  
895 of previous correct choices, and the number of previous incorrect choices,  
896 whose effects were also included. Lastly for task C, we included the slope  
897 of accuracy and response times over set sizes, as measures of the effect of  
898 set size on performance. For task B, we also included the difference between  
899 early (first third of trials) and late (last third) performance as a measure of  
900 learning. To avoid biases in the PCA toward any specific task, we included  
901 equal numbers of behavioral features for each task.

902 To facilitate the interpretation of PC2 and PC3, we normalized the load-  
903 ings (PCA weights) of each feature (behavioral and model parameter) with  
904 respect to PC1, flipping the loadings of all features in PC2 and PC3 that  
905 loaded negatively on PC1. This step ensured that the directions of factor  
906 loadings on PC2 and PC3 were interpretable in the same way for all features,  
907 irrespective of their role for task performance, and revealed the encoding of  
908 task contrasts.

#### 909 4.7. Ridge Regression

910 In sections 2.2.2 and 2.2.3, we use regularized, cross-validated Ridge re-  
911 gression to determine whether parameters captured overlapping variance,  
912 which would point to an overlap in cognitive processes. We used Ridge re-  
913 gression to avoid problems that would be caused by overfitting when using  
914 regular regression models. Ridge regression regularizes regression weight pa-  
915 rameters  $w$  based on their L2-norm. Regular regression identifies a vector  
916 of regression weights  $w$  that minimize the linear least squares  $\|y - wX\|_2^2$ .  
917 Here,  $\|a\|_2^2 = \sqrt{\sum_{a_i \in x} a_i^2}$  is the L2-norm of a vector  $a$ , vector  $y$  represents  
918 the outcome variable (in our case, a vector of parameters, one fitted to each  
919 participant), matrix  $X$  represents the predictor variables (in our case, either  
920 several behavioral features for each participant [2.2.2], or several parame-  
921 ters fitted to each participant 2.2.3]), and vector  $w$  represents the weights  
922 assigned to each feature in  $X$  (in our case, the weight assigned to each pre-  
923 dicting behavioral pattern or each predicting parameter).

924 When datasets are small compared to the number of predictors in a re-  
925 gression model, *exploding* regression weights  $w$  can lead to overfitting. Ridge  
926 regression avoids this issue by not only minimizing the linear least squares  
927 like regular regression, but also the L2 norm of weights  $w$ , i.e., by minimizing  
928  $\|y - wX\|_2^2 + \alpha * \|w\|_2^2$ . Parameter  $\alpha$  is a hyper-parameter of Ridge regression,



929 which needs to be chosen by the experimenter. To avoid bias in the selection  
930 of  $\alpha$ , we employed repeated cross-validated grid search. At each iteration of  
931 this procedure, we split the dataset into a predetermined number  $s \in [2, 3,$   
932  $\dots, 8]$  of equal-sized folds, and then fitted a Ridge regression to each fold, us-  
933 ing values of  $\alpha \in [0, 10, 30, 50, 100, 300, \dots, 10,000, 100,000, 1,000,000]$ . For  
934 each  $s$ , we determined the best value of  $\alpha$  based on cross-validation between  
935 folds, using the amount of explained variance,  $R^2$ , as the selection criterion.  
936 To avoid biases based on the random assignment of participants into folds,  
937 we repeated this procedure  $n = 100$  times for each value of  $\alpha$ . To avoid biases  
938 due to the number of folds, the entire process was repeated for each  $s$ , and  
939 the final value of  $s$  was selected based on  $R^2$ . We used the python package  
940 “scikit learn” [92] to implement the procedure.

941 We conducted three models per parameter to determine the relations be-  
942 tween parameters: predicting each parameter from all the parameters of each  
943 of the other two tasks (2 models); and predicting each parameter from all  
944 parameters of both other tasks combined (1 model; Fig. 4A). We conducted  
945 the same three models per parameter to determine the relations between pa-  
946 rameters and behaviors, predicting each parameter from behavioral features  
947 of the other tasks (Fig. 4A). In addition, we conducted a fourth model for  
948 behaviors, predicting each parameter from the behaviors of all three tasks  
949 combined, to assess the contributions of all behaviors to each parameter (Fig.  
950 4C). Meta-parameters  $s$  and  $\alpha$  were allowed to differ (and differed) between  
951 models. The final values of  $R^2$  (Fig. 4B and 4D) and the final regression  
952 weights  $w$  (Fig. 4A and 4C; Table 6) were determined by refitting the winning  
953 model.

## 954 5. Acknowledgments

955 We thank Ian Ballard, Mayank Agrawal, Gautam Agarwal, and Bas van  
956 Opheusden for helpful comments on this manuscript, and Catherine Hart-  
957 ley and other members of her lab for fruitful discussion. Numerous people  
958 contributed to this research: Amy Zou, Lance Kriegsfeld, Celia Ford, Jen-  
959 nifer Pfeifer, Megan Johnson, Vy Pham, Rachel Arsenault, Josephine Chris-  
960 ton, Shoshana Edelman, Lucy Eletel, Neta Gotlieb, Haley Keglovits, Julie  
961 Liu, Justin Morillo, Nithya Rajakumar, Nick Spence, Tanya Smith, Ben-  
962 jamin Tang, Talia Welte, and Lucy Whitmore. We are also grateful to our  
963 participants and their families. The work was funded by National Science  
964 Foundation SL-CN grant 1640885 to RD, AGEK, and LW.

965 **References**

- 966 [1] W. Schultz, P. Dayan, P. R. Montague, A Neural Substrate  
967 of Prediction and Reward, *Science* 275 (5306) (1997) 1593–1599.  
968 doi:10.1126/science.275.5306.1593.
- 969 [2] J. P. O’Doherty, P. Dayan, J. Schultz, R. Deichmann, K. Friston, R. J.  
970 Dolan, Dissociable Roles of Ventral and Dorsal Striatum in Instru-  
971 mental Conditioning, *Science* 304 (5669) (2004) 452–454, publisher:  
972 American Association for the Advancement of Science Section: Re-  
973 port. doi:10.1126/science.1094285.  
974 URL <https://science.sciencemag.org/content/304/5669/452>
- 975 [3] J. Gläscher, A. N. Hampton, J. P. O’Doherty, Determining a role for  
976 ventromedial prefrontal cortex in encoding action-based value signals  
977 during reward-related decision making, *Cerebral Cortex* (New York,  
978 N.Y.: 1991) 19 (2) (2009) 483–495. doi:10.1093/cercor/bhn098.
- 979 [4] J. Ribas Fernandes, A. Solway, C. Diuk, J. T. McGuire, A. G.  
980 Barto, Y. Niv, M. Botvinick, A Neural Signature of Hierar-  
981 chical Reinforcement Learning, *Neuron* 71 (2) (2011) 370–379.  
982 doi:10.1016/j.neuron.2011.05.042.
- 983 [5] N. Daw, S. Gershman, B. Seymour, P. Dayan, R. Dolan, Model-Based  
984 Influences on Humans’ Choices and Striatal Prediction Errors, *Neuron*  
985 69 (6) (2011) 1204–1215. doi:10.1016/j.neuron.2011.02.027.
- 986 [6] J. X. Wang, Z. Kurth-Nelson, D. Kumaran, D. Tirumala, H. Soyer,  
987 J. Z. Leibo, D. Hassabis, M. Botvinick, Prefrontal cortex as a meta-  
988 reinforcement learning system, *Nature Neuroscience* 21 (6) (2018) 860–  
989 868. doi:10.1038/s41593-018-0147-8.
- 990 [7] M. K. Eckstein, A. G. E. Collins, Computational evidence for hier-  
991 archically structured reinforcement learning in humans, *Proceedings*  
992 *of the National Academy of Sciences* 117 (47) (2020) 29381–29389.  
993 doi:10.1073/pnas.1912330117.  
994 URL <https://www.pnas.org/content/117/47/29381>
- 995 [8] M. Botvinick, Hierarchical reinforcement learning and decision  
996 making, *Current Opinion in Neurobiology* 22 (6) (2012) 956–962.

- 997 doi:10.1016/j.conb.2012.05.008.  
998 URL <http://linkinghub.elsevier.com/retrieve/pii/S0959438812000876>
- 999 [9] A. G. E. Collins, E. Koechlin, Reasoning, Learning, and Creativity:  
1000 Frontal Lobe Function and Human Decision-Making, *PLOS Biology*  
1001 10 (3) (2012) e1001293. doi:10.1371/journal.pbio.1001293.
- 1002 [10] D. M. Werchan, A. G. E. Collins, M. J. Frank, D. Amso, Role of  
1003 Prefrontal Cortex in Learning and Generalizing Hierarchical Rules  
1004 in 8-Month-Old Infants, *The Journal of Neuroscience* 36 (40) (2016)  
1005 10314–10322. doi:10.1523/JNEUROSCI.1351-16.2016.  
1006 URL <http://www.jneurosci.org/lookup/doi/10.1523/JNEUROSCI.1351-16.2016>
- 1007 [11] W. van den Bos, R. Bruckner, M. R. Nassar, R. Mata, B. Ep-  
1008 pinger, Computational neuroscience across the lifespan: Promises  
1009 and pitfalls, *Developmental Cognitive Neuroscience* (Oct. 2017).  
1010 doi:10.1016/j.dcn.2017.09.008.  
1011 URL <http://linkinghub.elsevier.com/retrieve/pii/S1878929317301068>
- 1012 [12] F. Bolenz, A. M. F. Reiter, B. Eppinger, Developmental Changes  
1013 in Learning: Computational Mechanisms and Social Influ-  
1014 ences, *Frontiers in Psychology* 8, publisher: Frontiers (2017).  
1015 doi:10.3389/fpsyg.2017.02048.  
1016 URL <https://www.frontiersin.org/articles/10.3389/fpsyg.2017.02048/full>
- 1017 [13] K. Nussenbaum, C. A. Hartley, Reinforcement learning across  
1018 development: What insights can we draw from a decade of re-  
1019 search?, *Developmental Cognitive Neuroscience* 40 (2019) 100733.  
1020 doi:10.1016/j.dcn.2019.100733.  
1021 URL <http://www.sciencedirect.com/science/article/pii/S1878929319303202>
- 1022 [14] Q. J. M. Huys, T. V. Maia, M. J. Frank, Computational psychiatry as a  
1023 bridge from neuroscience to clinical applications, *Nature neuroscience*  
1024 19 (3) (2016) 404–413. doi:10.1038/nn.4238.  
1025 URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5443409/>
- 1026 [15] R. A. Adams, Q. J. M. Huys, J. P. Roiser, Computational Psychiatry:  
1027 towards a mathematically informed understanding of mental illness,  
1028 *Journal of Neurology, Neurosurgery & Psychiatry* 87 (1) (2016) 53–  
1029 63, publisher: BMJ Publishing Group Ltd Section: Neuropsychiatry.

- 1030 doi:10.1136/jnnp-2015-310737.  
1031 URL <https://jnnp.bmj.com/content/87/1/53>
- 1032 [16] T. U. Hauser, G.-J. Will, M. Dubois, R. J. Dolan, Annual Re-  
1033 search Review: Developmental computational psychiatry, Jour-  
1034 nal of Child Psychology and Psychiatry 60 (4) (2019) 412–426.  
1035 doi:<https://doi.org/10.1111/jcpp.12964>.
- 1036 [17] W.-Y. Ahn, J. R. Busemeyer, Challenges and promises for trans-  
1037 lating computational tools into clinical practice, Current Opinion  
1038 in Behavioral Sciences 11 (2016) 1–7. doi:10.1016/j.cobeha.2016.02.001.  
1039 URL <https://www.sciencedirect.com/science/article/pii/S2352154616300237>
- 1040 [18] L. Deserno, R. Boehme, A. Heinz, F. Schlagenhauf, Reinforcement  
1041 Learning and Dopamine in Schizophrenia: Dimensions of Symptoms  
1042 or Specific Features of a Disease Group?, Frontiers in Psychiatry 4,  
1043 publisher: Frontiers (2013). doi:10.3389/fpsy.2013.00172.  
1044 URL <https://www.frontiersin.org/articles/10.3389/fpsy.2013.00172/full>
- 1045 [19] M. J. Frank, E. D. Claus, Anatomy of a decision: Striato-orbitofrontal  
1046 interactions in reinforcement learning, decision making, and reversal.,  
1047 Psychological Review 113 (2) (2006) 300–326. doi:10.1037/0033-  
1048 295X.113.2.300.  
1049 URL <http://doi.apa.org/getdoi.cfm?doi=10.1037/0033-295X.113.2.300>
- 1050 [20] Y. Niv, Reinforcement learning in the brain, Journal of Mathematical  
1051 Psychology 53 (3) (2009) 139–154.
- 1052 [21] D. Lee, H. Seo, M. W. Jung, Neural Basis of Reinforcement Learning  
1053 and Decision Making, Annual review of neuroscience 35 (2012) 287–  
1054 308. doi:10.1146/annurev-neuro-062111-150512.
- 1055 [22] J. P. O’Doherty, S. W. Lee, D. McNamee, The structure  
1056 of reinforcement-learning mechanisms in the human brain,  
1057 Current Opinion in Behavioral Sciences 1 (2015) 94–100.  
1058 doi:10.1016/j.cobeha.2014.10.004.
- 1059 [23] P. W. Glimcher, Understanding dopamine and reinforcement learning:  
1060 The dopamine reward prediction error hypothesis, Proceedings of the  
1061 National Academy of Sciences 108 (3) (2011) 15647–15654.

- 1062 [24] J. Garrison, B. Erdeniz, J. Done, Prediction error in reinforce-  
1063 ment learning: A meta-analysis of neuroimaging studies, *Neu-*  
1064 *roscience & Biobehavioral Reviews* 37 (7) (2013) 1297–1310.  
1065 doi:10.1016/j.neubiorev.2013.03.023.  
1066 URL <http://www.sciencedirect.com/science/article/pii/S0149763413000833>
- 1067 [25] P. Dayan, Y. Niv, Reinforcement learning: The Good, The Bad and  
1068 The Ugly, *Current Opinion in Neurobiology* 18 (2) (2008) 185–196.  
1069 doi:10.1016/j.conb.2008.08.003.  
1070 URL <https://linkinghub.elsevier.com/retrieve/pii/S0959438808000767>
- 1071 [26] R. S. Sutton, A. G. Barto, Reinforcement Learning: An Introduction,  
1072 2nd Edition, MIT Press, Cambridge, MA; London, England, 2017.
- 1073 [27] S. Palminteri, V. Wyart, E. Koechlin, The Importance of Falsification  
1074 in Computational Cognitive Modeling, *Trends in Cognitive Sciences*  
1075 21 (6) (2017) 425–433. doi:10.1016/j.tics.2017.03.011.  
1076 URL <https://linkinghub.elsevier.com/retrieve/pii/S1364661317300542>
- 1077 [28] M. K. Eckstein, L. Wilbrecht, A. G. E. Collins, What do Rein-  
1078 forcement Learning Models Measure? Interpreting Model Parameters  
1079 in Cognition and Neuroscience, *psyArxivType: article* (May 2021).  
1080 doi:10.31234/osf.io/e7kwx.  
1081 URL <https://psyarxiv.com/e7kwx/>
- 1082 [29] N. D. Daw, Trial-by-trial data analysis using computational models,  
1083 *Decision Making, Affect, and Learning: Attention and Performance*  
1084 *XXIII* (2011). doi:10.1093/acprof:oso/9780199600434.003.0001.
- 1085 [30] R. C. Wilson, A. G. Collins, Ten simple rules for the computational  
1086 modeling of behavioral data, *eLife* 8 (2019) e49547, publisher: eLife  
1087 Sciences Publications, Ltd. doi:10.7554/eLife.49547.  
1088 URL <https://doi.org/10.7554/eLife.49547>
- 1089 [31] O. Guest, A. E. Martin, How Computational Modeling Can Force The-  
1090 ory Building in Psychological Science, *Perspectives on Psychological*  
1091 *Science* (2021) 1745691620970585 Publisher: SAGE Publications Inc.  
1092 doi:10.1177/1745691620970585.  
1093 URL <https://doi.org/10.1177/1745691620970585>

- 1094 [32] G. Blohm, K. P. Kording, P. R. Schrater, A How-to-Model Guide  
1095 for Neuroscience, *eNeuro* 7 (1), publisher: Society for Neuro-  
1096 science Section: Research Article: Methods/New Tools (Jan. 2020).  
1097 doi:10.1523/ENEURO.0352-19.2019.  
1098 URL <https://www.eneuro.org/content/7/1/ENEURO.0352-19.2019>
- 1099 [33] S. J. Gershman, Empirical priors for reinforcement learning  
1100 models, *Journal of Mathematical Psychology* 71 (2016) 1–6.  
1101 doi:10.1016/j.jmp.2016.01.006.  
1102 URL <http://www.sciencedirect.com/science/article/pii/S0022249616000080>
- 1103 [34] T. Harada, Learning From Success or Failure? – Positiv-  
1104 ity Biases Revisited, *Frontiers in Psychology* 11 (Jul. 2020).  
1105 doi:10.3389/fpsyg.2020.01627.  
1106 URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7396482/>
- 1107 [35] S. L. Master, M. K. Eckstein, N. Gotlieb, R. Dahl, L. Wilbrecht,  
1108 A. G. E. Collins, Disentangling the systems contributing to changes  
1109 in learning during adolescence, *Developmental Cognitive Neuroscience*  
1110 41 (2020) 100732. doi:10.1016/j.dcn.2019.100732.  
1111 URL <http://www.sciencedirect.com/science/article/pii/S1878929319303196>
- 1112 [36] Z. A. Yaple, R. Yu, Fractionating adaptive learning: A meta-analysis  
1113 of the reversal learning paradigm, *Neuroscience & Biobehavioral*  
1114 *Reviews* 102 (2019) 85–94. doi:10.1016/j.neubiorev.2019.04.006.  
1115 URL <http://www.sciencedirect.com/science/article/pii/S0149763418308996>
- 1116 [37] X. Liu, J. Hairston, M. Schrier, J. Fan, Common and distinct networks  
1117 underlying reward valence and processing stages: A meta-analysis of  
1118 functional neuroimaging studies, *Neuroscience and Biobehavioral Re-*  
1119 *views* 35 (5) (2011) 1219–1236. doi:10.1016/j.neubiorev.2010.12.012.  
1120 URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3395003/>
- 1121 [38] W. Schultz, A. Dickinson, Neuronal Coding of Prediction Er-  
1122 rors, *Annual Review of Neuroscience* 23 (1) (2000) 473–500.  
1123 doi:10.1146/annurev.neuro.23.1.473.  
1124 URL <http://www.annualreviews.org/doi/10.1146/annurev.neuro.23.1.473>
- 1125 [39] A. G. E. Collins, J. K. Brown, J. M. Gold, J. A. Waltz, M. J.  
1126 Frank, Working Memory Contributions to Reinforcement Learning Im-  
1127 pairments in Schizophrenia, *Journal of Neuroscience* 34 (41) (2014)

1128 13747–13756, publisher: Society for Neuroscience Section: Articles.  
1129 doi:10.1523/JNEUROSCI.0989-14.2014.  
1130 URL <https://www.jneurosci.org/content/34/41/13747>

[40] A. H. Javadi, D. H. K. Schmidt, M. N. Smolka, Adolescents adapt more slowly than adults to varying reward contingencies, *Journal of Cognitive Neuroscience* 26 (12) (2014) 2670–2681. doi:10.1162/jocn\_a00677.

[41] S.-J. Blakemore, T. W. Robbins, Decision-making in the adolescent brain, *Nature Neuroscience* 15 (9) (2012) 1184–1191, number: 9 Publisher: Nature Publishing Group. doi:10.1038/nn.3177.  
1132  
1133  
1134 URL <http://www.nature.com/articles/nn.3177>

[42] S. DePasque, A. Galván, Frontostriatal development and probabilistic reinforcement learning during adolescence, *Neurobiology of Learning and Memory* 143 (2017) 1–7. doi:10.1016/j.nlm.2017.04.009.  
1136  
1137  
1138 URL <http://www.sciencedirect.com/science/article/pii/S107474271730062X>

[43] A. Mohebi, J. R. Pettibone, A. A. Hamid, J.-M. T. Wong, L. T. Vinson, T. Patriarchi, L. Tian, R. T. Kennedy, J. D. Berke, Dissociable dopamine dynamics for learning and motivation, *Nature* 570 (7759) (2019) 65–70, number: 7759 Publisher: Nature Publishing Group. doi:10.1038/s41586-019-1235-y.  
1140  
1141  
1142  
1143 URL <https://www.nature.com/articles/s41586-019-1235-y>

[44] W. R. Uttal, On some two-way barriers between models and mechanisms, *Perception & Psychophysics* 48 (2) (1990) 188–203. doi:10.3758/BF03207086.  
1145  
1146 URL <https://doi.org/10.3758/BF03207086>

[45] B. Webb, Can robots make good models of biological behaviour?, *Behavioral and Brain Sciences* 24 (6) (2001) 1033–1050, publisher: Cambridge University Press. doi:10.1017/S0140525X01000127.  
1148  
1149  
1150 URL <http://www.cambridge.org/core/journals/behavioral-and-brain-sciences/article/>

[46] D. J. Navarro, Between the Devil and the Deep Blue Sea: Tensions Between Scientific Judgement and Statistical Model Selection, *Computational Brain & Behavior* 2 (1) (2019) 28–34. doi:10.1007/s42113-018-0019-z.  
1152  
1153  
1154 URL <https://doi.org/10.1007/s42113-018-0019-z>

[47] T. Yarkoni, The generalizability crisis, *The Behavioral and brain sciences* Publisher: Behav Brain Sci (Dec. 2020). doi:10.1017/S0140525X20001685.  
1156  
1157 URL <https://pubmed.ncbi.nlm.nih.gov/33342451/>



- [48] V. M. Brown, J. Chen, C. M. Gillan, R. B. Price, Improving the Reliability  
1159 of Computational Analyses: Model-Based Planning and Its Relationship  
1160 With Compulsivity, *Biological Psychiatry: Cognitive Neuroscience and*  
1161 *Neuroimaging* 5 (6) (2020) 601–609. doi:10.1016/j.bpsc.2019.12.019.  
1162 URL <https://www.sciencedirect.com/science/article/pii/S2451902220300161>
- [49] M. K. Eckstein, S. L. Master, R. E. Dahl, L. Wilbrecht, A. G. E. Collins,  
1164 Understanding the Unique Advantage of Adolescents in Stochastic, Volatile  
1165 Environments: Combining Reinforcement Learning and Bayesian Inference,  
1166 *bioRxiv* (2020) 2020.07.04.187971 Publisher: Cold Spring Harbor Laboratory  
1167 Section: New Results. doi:10.1101/2020.07.04.187971.  
1168 URL <https://www.biorxiv.org/content/10.1101/2020.07.04.187971v1>
- [50] L. Xia, S. Master, M. Eckstein, L. Wilbrecht, A. G. E. Collins, Learning under  
1170 uncertainty changes during adolescence, in: *Proceedings of the Cognitive*  
1171 *Science Society*, 2020.
- [51] M. D. Lee, How cognitive modeling can benefit from hierarchical  
1173 Bayesian models, *Journal of Mathematical Psychology* 55 (1) (2011)  
1174 1–7. doi:10.1016/j.jmp.2010.08.013.  
1175 URL <https://linkinghub.elsevier.com/retrieve/pii/S0022249610001148>
- [52] A. G. E. Collins, M. J. Frank, How much of reinforcement learning is work-  
1177 ing memory, not reinforcement learning? A behavioral, computational, and  
1178 neurogenetic analysis: Working memory in reinforcement learning, *Euro-*  
1179 *pean Journal of Neuroscience* 35 (7) (2012) 1024–1035. doi:10.1111/j.1460-  
1180 9568.2011.07980.x.
- [53] L. H. Somerville, S. F. Sasse, M. C. Garrad, A. T. Drysdale, N. Abi Akar,  
1182 C. Insel, R. C. Wilson, Charting the expansion of strategic exploratory behav-  
1183 ior during adolescence, *Journal of Experimental Psychology: General* 146 (2)  
1184 (2017) 155–164, place: US Publisher: American Psychological Association.  
1185 doi:10.1037/xge0000250.
- [54] A. Gopnik, Childhood as a solution to explore–exploit tensions, *Philosoph-*  
1187 *ical Transactions of the Royal Society B: Biological Sciences* 375 (1803)  
1188 (2020) 20190502, publisher: Royal Society. doi:10.1098/rstb.2019.0502.  
1189 URL <https://royalsocietypublishing.org/doi/10.1098/rstb.2019.0502>

- [55] T. E. J. Behrens, M. W. Woolrich, M. E. Walton, M. F. S. Rushworth,  
1191 Learning the value of information in an uncertain world, *Nature Neuroscience*  
1192 10 (9) (2007) 1214–1221. doi:10.1038/nn1954.  
1193 URL <https://www.nature.com/articles/nn1954>
- [56] K. C. Berridge, The debate over dopamine’s role in reward: the  
1195 case for incentive salience, *Psychopharmacology* 191 (3) (2007) 391–431.  
1196 doi:10.1007/s00213-006-0578-x.  
1197 URL <https://doi.org/10.1007/s00213-006-0578-x>
- [57] A. J. Yu, P. Dayan, Uncertainty, Neuromodulation, and Attention, *Neuron*  
1199 46 (4) (2005) 681–692. doi:10.1016/j.neuron.2005.04.026.  
1200 URL <http://www.sciencedirect.com/science/article/pii/S0896627305003624>
- [58] S. Bouret, S. J. Sara, Network reset: a simplified overarching theory of locus  
1202 coeruleus noradrenaline function, *Trends in Neurosciences* 28 (11) (2005)  
1203 574–582. doi:10.1016/j.tins.2005.09.002.
- [59] S. J. Gershman, Dopamine, Inference, and Uncertainty, *Neural Computation*  
29 (12) (2017) 3311–3326. doi:10.1162/neco\_a01023.  
URL [http://www.mitpressjournals.org/doi/abs/10.1162/neco\\_a01023](http://www.mitpressjournals.org/doi/abs/10.1162/neco_a01023). *K. Starkweather, S.*  
1204
- [60] S. J. Gershman, N. Uchida, Believing in dopamine, *Nature Reviews Neuro-*  
1206 *science* 20 (11) (2019) 703–714, number: 11 Publisher: Nature Publishing  
1207 Group. doi:10.1038/s41583-019-0220-7.  
1208 URL <https://www.nature.com/articles/s41583-019-0220-7>
- [62] K. Katahira, The statistical structures of reinforcement learning with  
1210 asymmetric value updates, *Journal of Mathematical Psychology* 87 (2018)  
1211 31–45. doi:10.1016/j.jmp.2018.09.002.  
1212 URL <http://www.sciencedirect.com/science/article/pii/S0022249617302407>
- [63] M. Sugawara, K. Katahira, Dissociation between asymmetric value up-  
1214 dating and perseverance in human reinforcement learning, *Scientific Re-*  
1215 *ports* 11 (1) (2021) 3574, number: 1 Publisher: Nature Publishing Group.  
1216 doi:10.1038/s41598-020-80593-7.  
1217 URL <https://www.nature.com/articles/s41598-020-80593-7>

- [64] C. Diuk, A. Schapiro, N. Córdoba, J. Ribas-Fernandes, Y. Niv, M. Botvinick, Divide and Conquer: Hierarchical Reinforcement Learning and Task Decomposition in Humans, in: Computational and Robotic Models of the Hierarchical Organization of Behavior, Springer, Berlin, Heidelberg, 2013, pp. 271–291. doi:10.1007/978-3-642-39875-9<sub>1</sub>2.
- [65] A. G. E. Collins, The Tortoise and the Hare: Interactions between Reinforcement Learning and Working Memory, *Journal of Cognitive Neuroscience* 30 (10) (2018) 1422–1432. doi:10.1162/jocn<sub>a</sub>01238.
- [66] L.-H. Tai, A. M. Lee, N. Benavidez, A. Bonci, L. Wilbrecht, Transient stimulation of distinct subpopulations of striatal neurons mimics changes in action value, *Nature Neuroscience* 15 (9) (2012) 1281–1289. doi:10.1038/nn.3188.
- [67] M. Donoso, A. G. E. Collins, E. Koechlin, Foundations of human reasoning in the prefrontal cortex, *Science* 344 (6191) (2014) 1481–1486. doi:10.1126/science.1252254.  
URL <http://www.sciencemag.org/cgi/doi/10.1126/science.1252254>
- [68] A. M. Bornstein, K. A. Norman, Reinstated episodic context guides sampling-based decisions for reward, *Nature Neuroscience* 20 (7) (2017) 997–1003. doi:10.1038/nn.4573.  
URL <https://www.nature.com/articles/nn.4573>
- [69] O. M. Vikbladh, M. R. Meager, J. King, K. Blackmon, O. Devinsky, D. Shohamy, N. Burgess, N. D. Daw, Hippocampal Contributions to Model-Based Planning and Spatial Memory, *Neuron* 102 (3) (2019) 683–693.e4. doi:10.1016/j.neuron.2019.02.014.  
URL <https://www.sciencedirect.com/science/article/pii/S0896627319301230>
- [70] M. E. van der Schaaf, E. Warmerdam, E. A. Crone, R. Cools, Distinct linear and non-linear trajectories of reward and punishment reversal learning during development: relevance for dopamine’s role in adolescent decision making, *Developmental Cognitive Neuroscience* 1 (4) (2011) 578–590. doi:10.1016/j.dcn.2011.06.007.
- [71] N. Sendhilnathan, M. Semework, M. E. Goldberg, A. E. Ipata, Neural Correlates of Reinforcement Learning in Mid-lateral Cerebellum, *Neuron* 106 (1) (2020) 188–198.e5. doi:10.1016/j.neuron.2019.12.032.

- [72] S. D. McDougle, A. G. E. Collins, Modeling the influence of working memory, reinforcement, and action uncertainty on reaction time and choice during instrumental learning, *Psychonomic Bulletin & Review* 28 (1) (2021) 20–39. doi:10.3758/s13423-020-01774-z. URL <https://doi.org/10.3758/s13423-020-01774-z>
- [73] A. Radulescu, Y. Niv, I. Ballard, Holistic Reinforcement Learning: The Role of Structure and Attention, *Trends in Cognitive Sciences* 23 (4) (2019) 278–292. doi:10.1016/j.tics.2019.01.010. URL <https://www.sciencedirect.com/science/article/pii/S1364661319300361>
- [74] A. Kononov, I. Krajbich, Neurocomputational Dynamics of Sequence Learning, *Neuron* 98 (6) (2018) 1282–1293.e4. doi:10.1016/j.neuron.2018.05.013. URL <http://www.sciencedirect.com/science/article/pii/S0896627318303854>
- [75] W. Kool, F. A. Cushman, S. J. Gershman, When Does Model-Based Control Pay Off?, *PLOS Computational Biology* 12 (8) (2016) e1005090, publisher: Public Library of Science. doi:10.1371/journal.pcbi.1005090. URL <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1005090>
- [76] C. F. d. Silva, T. A. Hare, Humans are primarily model-based learners in the two-stage task, *bioRxiv* (2020) 682922Publisher: Cold Spring Harbor Laboratory Section: New Results. doi:10.1101/682922. URL <https://www.biorxiv.org/content/10.1101/682922v4>
- [77] I. W. Eisenberg, P. G. Bissett, A. Zeynep Enkavi, J. Li, D. P. MacKinnon, L. A. Marsch, R. A. Poldrack, Uncovering the structure of self-regulation through data-driven ontology discovery, *Nature Communications* 10 (1) (2019) 1–13. doi:10.1038/s41467-019-10301-1. URL <https://www.nature.com/articles/s41467-019-10301-1>
- [78] W. C. Lin, K. Delevich, L. Wilbrecht, A role for adaptive developmental plasticity in learning and decision making, *Current Opinion in Behavioral Sciences* 36 (2020) 48–54. doi:10.1016/j.cobeha.2020.07.010. URL <http://www.sciencedirect.com/science/article/pii/S2352154620301121>
- [79] J. Davidow, K. Foerde, A. Galvan, D. Shohamy, An Upside to Reward Sensitivity: The Hippocampus Supports Enhanced Reinforcement Learning in

- 1274 Adolescence, *Neuron* 92 (1) (2016) 93–99. doi:10.1016/j.neuron.2016.08.031.  
1275 URL <http://linkinghub.elsevier.com/retrieve/pii/S0896627316305244>
- [80] C. Johnson, L. Wilbrecht, Juvenile mice show greater flexibility in multiple  
1277 choice reversal learning than adults, *Developmental Cognitive Neuroscience*  
1278 1 (4) (2011) 540–551. doi:10.1016/j.dcn.2011.05.008.  
1279 URL <http://linkinghub.elsevier.com/retrieve/pii/S1878929311000533>
- [81] A. C. Petersen, L. Crockett, M. Richards, A. Boxer, A self-report measure  
1281 of pubertal status: Reliability, validity, and initial norms, *Journal of Youth  
1282 and Adolescence* 17 (2) (1988) 117–133. doi:10.1007/BF01537962.  
1283 URL <https://doi.org/10.1007/BF01537962>
- [82] A. Christakou, S. J. Gershman, Y. Niv, A. Simmons, M. Brammer, K. Rubia,  
Neural and psychological maturation of decision-making in adolescence and  
young adulthood, *Journal of Cognitive Neuroscience* 25 (11) (2013) 1807–  
1823. doi:10.1162/jocn\_a00447.
- [83] W. van den Bos, M. X. Cohen, T. Kahnt, E. A. Crone, Striatum–Medial  
1285 Prefrontal Cortex Connectivity Predicts Developmental Changes in  
1286 Reinforcement Learning, *Cerebral Cortex* 22 (6) (2012) 1247–1255.  
1287 doi:10.1093/cercor/bhr198.  
1288 URL <https://academic.oup.com/cercor/article/22/6/1247/307075>
- [84] M. J. Frank, L. C. Seeberger, R. C. O’Reilly, By Carrot or by Stick: Cognitive  
1290 Reinforcement Learning in Parkinsonism, *Science* 306 (5703) (2004) 1940–  
1291 1943. doi:10.1126/science.1102941.
- [85] R. D. Cazé, M. A. A. van der Meer, Adaptive properties of differential learn-  
1293 ing rates for positive and negative outcomes, *Biological Cybernetics* 107 (6)  
1294 (2013) 711–719. doi:10.1007/s00422-013-0571-5.
- [86] S. Palminteri, E. J. Kilford, G. Coricelli, S.-J. Blakemore, The Computational  
1296 Development of Reinforcement Learning during Adolescence, *PLoS Computa-  
1297 tional Biology* 12 (6) (Jun. 2016). doi:10.1371/journal.pcbi.1004953.  
1298 URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4920542/>
- [87] G. Lefebvre, M. Lebreton, F. Meyniel, S. Bourgeois-Gironde, S. Palminteri,  
1300 Behavioural and neural characterization of optimistic reinforcement learning,  
1301 *Nature Human Behaviour* 1 (4) (2017) 0067. doi:10.1038/s41562-017-0067.  
1302 URL <http://www.nature.com/articles/s41562-017-0067>

- [88] W. Dabney, Z. Kurth-Nelson, N. Uchida, C. K. Starkweather, D. Hassabis, R. Munos, M. Botvinick, A distributional code for value in dopamine-based reinforcement learning, *Nature* 577 (7792) (2020) 671–675. doi:10.1038/s41586-019-1924-6.  
URL <http://www.nature.com/articles/s41586-019-1924-6>
- [89] A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, D. B. Rubin, *Bayesian Data Analysis*, 3rd Edition, Chapman and Hall/CRC, Boca Raton, 2013.
- [90] K. Katahira, How hierarchical models improve point estimates of model parameters at the individual level, *Journal of Mathematical Psychology* 73 (2016) 37–58. doi:10.1016/j.jmp.2016.03.007.
- [91] S. Watanabe, A Widely Applicable Bayesian Information Criterion, *Journal of Machine Learning Research* 14 (Mar) (2013) 867–897.  
URL <http://www.jmlr.org/papers/v14/watanabe13a.html>
- [92] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, Duchesnay, Scikit-learn: Machine Learning in Python, *Journal of Machine Learning Research* 12 (85) (2011) 2825–2830.  
URL <http://jmlr.org/papers/v12/pedregosa11a.html>
- [93] D. A. Peterson, C. Elliott, D. D. Song, S. Makeig, T. J. Sejnowski, H. Poizner, Probabilistic reversal learning is impaired in Parkinson’s disease, *Neuroscience* 163 (4) (2009) 1092–1101. doi:10.1016/j.neuroscience.2009.07.033.  
URL <http://www.sciencedirect.com/science/article/pii/S0306452209012068>
- [94] R. Swainson, R. D. Rogers, B. J. Sahakian, B. A. Summers, C. E. Polkey, T. W. Robbins, Probabilistic learning and reversal deficits in patients with Parkinson’s disease or frontal or temporal lobe lesions: possible adverse effects of dopaminergic medication, *Neuropsychologia* 38 (5) (2000) 596–612. doi:10.1016/S0028-3932(99)00103-7.  
URL <http://www.sciencedirect.com/science/article/pii/S0028393299001037>
- [95] J. A. Waltz, J. M. Gold, Probabilistic reversal learning impairments in schizophrenia: Further evidence of orbitofrontal dysfunction, *Schizophrenia*

- 1335 Research 93 (1) (2007) 296–303. doi:10.1016/j.schres.2007.03.010.  
1336 URL <http://www.sciencedirect.com/science/article/pii/S092099640700120X>
- [96] D. P. Dickstein, E. C. Finger, M. A. Brotman, B. A. Rich, D. S. Pine,  
1338 J. R. Blair, E. Leibenluft, Impaired probabilistic reversal learning in youths  
1339 with mood and anxiety disorders, *Psychological Medicine* 40 (7) (2010)  
1340 1089–1100. doi:10.1017/S0033291709991462.  
1341 URL <http://www.cambridge.org/core/journals/psychological-medicine/article/impair>
- [97] R. Cools, L. Clark, A. M. Owen, T. W. Robbins, Defining the Neural Mech-  
1342 anisms of Probabilistic Reversal Learning Using Event-Related Functional  
1343 Magnetic Resonance Imaging, *Journal of Neuroscience* 22 (11) (2002) 4563–  
1344 4567. doi:10.1523/JNEUROSCI.22-11-04563.2002.  
1345 URL <https://www.jneurosci.org/content/22/11/4563>
- [98] R. Cools, M. J. Frank, S. E. Gibbs, A. Miyakawa, W. Jagust, M. D’Esposito,  
1346 Striatal Dopamine Predicts Outcome-Specific Reversal Learning and Its Sen-  
1347 sitivity to Dopaminergic Drug Administration, *Journal of Neuroscience* 29 (5)  
1348 (2009) 1538–1543. doi:10.1523/JNEUROSCI.4467-08.2009.  
1349 URL <https://www.jneurosci.org/content/29/5/1538>
- [99] F. Lourenco, B. Casey, Adjusting behavior to changing environmental  
1350 demands with development, *Neuroscience & Biobehavioral Reviews* 37 (9)  
1351 (2013) 2233–2242. doi:10.1016/j.neubiorev.2013.03.003.  
1352 URL <https://linkinghub.elsevier.com/retrieve/pii/S0149763413000638>
- [100] A. Izquierdo, J. L. Brigman, A. K. Radke, P. H. Rudebeck, A. Holmes, The  
1353 neural basis of reversal learning: An updated perspective, *Neuroscience* 345  
1354 (2017) 12–26. doi:10.1016/j.neuroscience.2016.03.021.  
1355 URL <http://www.sciencedirect.com/science/article/pii/S030645221600244X>
- [101] A. G. E. Collins, B. Ciullo, M. J. Frank, D. Badre, Working Memory Load  
1356 Strengthens Reward Prediction Errors, *The Journal of Neuroscience* 37 (16)  
1357 (2017) 4332–4342. doi:10.1523/JNEUROSCI.2700-16.2017.  
1358 URL <http://www.jneurosci.org/lookup/doi/10.1523/JNEUROSCI.2700-16.2017>
- [102] A. G. E. Collins, M. A. Albrecht, J. A. Waltz, J. M. Gold, M. J.  
1359 Frank, Interactions Among Working Memory, Reinforcement Learning,  
1360 and Effort in Value-Based Choice: A New Paradigm and Selective  
1361 Deficits in Schizophrenia, *Biological Psychiatry* 82 (6) (2017) 431–439.  
1362  
1363



1368 doi:10.1016/j.biopsych.2017.05.017.

1369 URL <http://www.sciencedirect.com/science/article/pii/S0006322317316190>