

1 Understanding the Unique Advantage of Adolescents in
2 Stochastic, Volatile Environments: Combining
3 Reinforcement Learning and Bayesian Inference

4 **Maria K. Eckstein¹, Sarah L. Master¹, Ronald E. Dahl², Linda**
5 **Wilbrecht^{1,3}, and Anne G.E. Collins¹**

6 ¹Department of Psychology, 2121 Berkeley Way West

7 ²Institute of Human Development, 2121 Berkeley Way West

8 ³Helen Wills Neuroscience Institute, 175 Li Ka Shing Center
9 Berkeley, California 94720 USA
10

11 **1. Abstract**

12 During adolescence, youth venture out, explore the wider world, and are
13 challenged to learn how to navigate novel and uncertain environments. We
14 investigated whether adolescents are uniquely adapted to this transition, com-
15 pared to younger children and adults. In a stochastic, volatile learning task
16 with a sample of 291 participants aged 8-30, we found that adolescents 13-15
17 years old outperformed both younger and older participants. We developed
18 two independent cognitive models, and used hierarchical Bayesian model fit-
19 ting to assess developmental changes in underlying cognitive mechanisms.
20 Choice parameters in both models improved monotonously. By contrast, up-
21 date parameters peaked closest to optimal values in 13-15 year-olds. Com-
22 bining both models using principal component analysis yielded new insights,
23 revealing that three components contributed to the early to mid-adolescent
24 performance peak. This research highlights early to mid-adolescence as a
25 neurodevelopmental window that may be more optimal for behavioral ad-
26 justment in volatile and uncertain environments. It also shows how detailed
27 insights can be gleaned by combining cognitive models.

28 **Keywords:** reinforcement learning, Bayesian inference, computational mod-
29 eling, development, volatility

2. Introduction

In mammals and other species with parental care, there is typically an adolescent stage of development in which the young are no longer supported by parental care but are not yet adult. This adolescent period can be identified in many species across the animal kingdom (Natterson-Horowitz and Bowers, 2019) and is increasingly viewed as a critical epoch of development in which organisms explore the world, make critical decisions, and learn about important features of their environment (DePasque and Galván, 2017; Laube et al., 2020; Piekarski, Johnson, et al., 2017; Steinberg, 2005). All of these behaviors require learning and decision making that will likely have critical short and long-term impact on survival of the organism (Frankenhuis and Walasek, 2020). In humans, and likely many other species, the transition to independence almost always involves environmental changes and increased exposure to stochastic, uncertain outcomes. It is therefore possible that adolescent brains and cognitive capabilities are specifically adapted to succeed in such situations (Dahl et al., 2018; Davidow et al., 2016; Johnson and Wilbrecht, 2011; Lourenco and Casey, 2013; Sercombe, 2014).

To test this idea, we compared the behaviors of 291 participants, including 191 children and adolescents aged 8-17, and 112 adults (55 adults from the community, aged 25-30; 57 university undergraduates, aged 18-28; suppl. Fig. 6), on a task with volatile structure and stochastic outcomes (Fig. 1A, B). The goal of the task was to collect rewards, which were hidden in one of two locations (Fig. 1A). Which location was rewarding changed unpredictably several times ("task switch"), and the rewarded location provided rewards only 75 percent of the time (Fig. 1A). The task's main challenge lay in discriminating chance outcomes during stable task periods from task switches, and respond appropriately to each. It required the integration of stochastic feedback and the adaptation to a volatile environment, and thus theoretically mirrored the challenges of the adolescent period. We therefore hypothesized that adolescents would outperform both younger and older participants. Our data supported this hypothesis.

We used computational modeling to understand the cognitive processes that underlie adolescents' superior performance, as well as the strategies employed by younger children and older teenagers and adults. A variety of algorithms have been used to model human cognition, including Reinforcement learning (RL) and Bayesian inference (BI). The basic idea of RL is that choice options have "values" (their expected long-term cumulative reward). The goal of RL—maximizing long-term outcomes—can therefore be achieved by selecting options according to their values. The core of RL lies in approximating values accurately and efficiently,

67 which can be achieved by performing small incremental updates every time an
68 outcome is observed. This incremental procedure avoids overemphasizing any
69 single outcome and allows RL to treat stochastic outcomes appropriately. The
70 size of the increment captures the integration time scale: the emphasis given to
71 recent vs. less recent outcomes. In volatile environments, RL adjusts to abrupt
72 changes by gradually unlearning and relearning values.

73 RL frames our task as a *learning* problem: Participants continuously learn
74 and adjust the value of each choice option based on trial-by-trial feedback (Fig.
75 3A, left). The same learning process occurs during stable periods and after task
76 switches, without an explicit concept of switching: Behavioral change arises when
77 enough updates have occurred for the values of one option to dip below the other.
78 Basic RL algorithms are suboptimal in volatile and structured environments like
79 ours, but can be augmented for more efficient performance (see Methods). In all
80 cases, RL models make the fundamental assumption that humans solve challenges
81 through continuous, value-based learning.

82 Computational modeling studies often select one type of cognitive model (e.g.,
83 RL), and compare different variants of this type to find the best-fitting one, which
84 is then interpreted as the cognitive process employed by participants. This ap-
85 proach cannot rule out whether a model of a different type (e.g., BI) would fit
86 the data better altogether. The issue can be mitigated by verifying that the chosen
87 model reproduces human behavior qualitatively and quantitatively, and a different
88 model type is unlikely to lead to even better results (Palminteri et al., 2017). How-
89 ever, a more troublesome concern is that different types of models frame behavior
90 in terms of different cognitive processes, and one framing can be more informa-
91 tive, more interpretable, or summarize behavior in a more meaningful way than
92 another. This problem is more difficult to solve because conceptual model fit is
93 hard to quantify.

94 To address these concerns, we fitted two families of models to the current
95 task, RL and BI. BI models combine "prior" knowledge with new observations to
96 arrive at "posterior" conclusions about unobservable features of the environment
97 ("hidden states"; Perfors et al., 2011; Sarkka, 2013). BI models therefore employ
98 a "predictive model", which specifies how likely different observations are to arise
99 from different hidden states ("likelihood"). The BI inference cycle—combining
100 prior and likelihood to get a posterior—can continue infinitely, using each step's
101 posterior as the prior for the subsequent step. Bayesian models deal well with
102 stochastic outcomes because extreme likelihoods are balanced by stable priors.
103 Environmental volatility is modeled explicitly as a change in hidden state.

104 BI models frame our task as an *inference* problem: Participants know that

105 the task has two hidden states ("Left choice is correct" and "Right choice is cor-
106 rect"; Fig. 3A, right), and use trial-by-trial outcomes to determine which state is
107 more likely. Having inferred the state, the appropriate action (left or right) can
108 be selected. In other words, participants entertain a mental model of the task,
109 which specifies how likely each outcome (reward, no reward) is in each hidden
110 state, and how likely state transitions occur. In summary, whereas RL claims that
111 participants adapt to task switches by continuously relearning choice values, BI
112 claims that they represent state transitions explicitly, changing their behavior after
113 detecting a switch.

114 We used the BI model to assess how participants' mental models developed
115 with age. We hypothesized that adolescents' models would be better tuned for
116 volatile and stochastic environments than children's and adults'. Because the BI
117 model employed rational, Bayes-optimal behavior, it also allowed us to evaluate
118 whether and how participants deviated from it: We hypothesized that adolescents
119 would use the most accurate mental models. In addition, both RL and BI models
120 contained parameters that controlled choice: decision noise and persistence. We
121 expected both to decrease monotonously with age, as has been consistently ob-
122 served (e.g., Master et al., 2019; for review, see Nussenbaum and Hartley, 2019).
123 For RL learning-rate parameters, which control integration time scales, we did
124 not have a priori predictions because past studies differed in experimental con-
125 text (Davidow et al., 2016; Master et al., 2019) and provided conflicting results
126 (Nussenbaum and Hartley, 2019).

127 Model-agnostic analyses revealed that adolescents (13-15 years) outperformed
128 younger and older participants in several measures of task performance, as pre-
129 dicted. We used state-of-the-art hierarchical Bayesian methods to fit RL and BI
130 models to participant behavior, assessing age changes directly and in a statisti-
131 cally unbiased way (Methods; Katahira, 2016; M. D. Lee, 2011; van den Bos et
132 al., 2017). Both models qualitatively captured participants' behavior, and choice-
133 related parameters showed the expected age trajectories. The BI model confirmed
134 the unique tuning of adolescents' (13-15) mental model to the task, and the RL
135 model revealed complex developmental trajectories of learning rates. Going be-
136 yond individual models, we then used Principal Component Analysis (PCA) to
137 expose the dimensions of largest variance in the shared parameter space. Variance
138 between participants was captured in just four dimensions, three of which showed
139 marked and separable developmental changes.

140 3. Results

141 3.1. Task

142 After completing a child-friendly tutorial (Methods), participants performed
143 the following task: On each trial, two identical green boxes appeared on the
144 screen. Participants chose one, and either received a reward (gold coin) or not
145 (empty box; Fig. 1A). One box was rewarded in 75% of the trials on which it was
146 chosen, whereas the other was never rewarded (*stochastic* aspect). After a vari-
147 able number of trials, an unsignaled switch occurred, after which the opposite box
148 was rewarding. Several unpredictable switches occurred over 120 trials (*volatility*
149 aspect; Fig. 1B). Participants' goal was to collect as many gold coins as possible.
150 More task details are provided in the Methods.

151 3.2. Task Behavior

152 Participants gradually adjusted their behavior after task switches, and on av-
153 erage started selecting the correct action about 3 trials after a switch, reaching
154 asymptotic performance thereafter (Fig. 1C). Participants almost always repeated
155 actions ("stayed") after receiving positive outcomes ("- +" and "+ +"), and often
156 switched actions after receiving two negative outcomes ("- -"). Behavior was am-
157 bivalent after receiving a positive followed by a negative outcome ("+ -"), i.e., on
158 "potential" switch trials (Fig. 1D; for age differences, see suppl. Fig. 15).

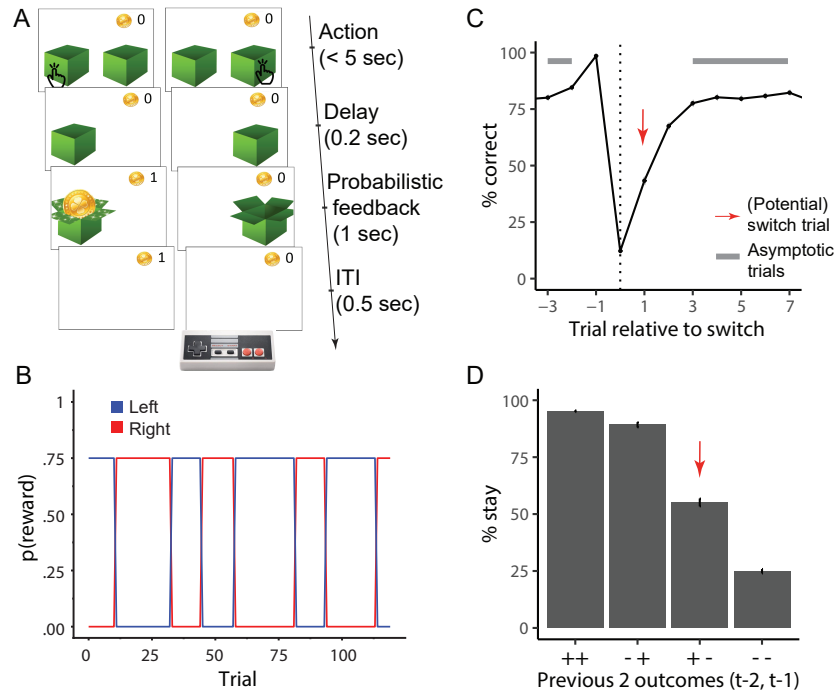


Figure 1: (A) Task design. On each trial, participants chose one of two boxes, using the two red buttons of the shown game controller. The chosen box either revealed a gold coin (left) or was empty (right). The probability of coin reward was 75% on the rewarded side, and 0% on the non-rewarded side. (B) The rewarded side changed multiple times, according to unpredictable task switches. (C) Average human performance and standard errors, aligned to "true" task switches (dotted line; trial 0). Switches only occurred after rewarded trials (Methods), resulting in performance of 100% on trial -1. The red arrow shows the "switch trial", grey bars show trials of asymptotic performance. (D) Average probability of repeating a previous action ("stay"), as a function of the two previous outcomes ($t-2$, $t-1$) for this action ("+": reward; "-": no reward). Error bars indicate between-participant standard errors. Red arrow highlights "potential switch trials", i.e., when a rewarded trial is followed by a non-rewarded one, which—from participants' perspective—is consistent with a task switch.

159 Focusing on age differences, adolescents 13-15 outperformed younger groups
 160 age 8-13 and adults (18-30) on several measures of performance (Fig. 2, suppl.
 161 Fig. 13, Fig. 3C-F). We tested age effects statistically with (logistic) mixed-effects
 162 regression (Methods). All measures of performance showed positive linear effects
 163 of age, indicating improved performance with age, as well as negative quadratic
 164 effects, consistent with a U-shaped relationship where adolescents 13-15 perform
 165 the task more accurately than both younger or older participants (Table 1).

Table 1: Statistics of mixed-effects regression models predicting performance measures from sex (male, female), age (years and months; "lin."), and squared age ("qua."). Overall accuracy, stay after potential (pot.) switch, and asymptotic performance were modeled using logistic regression, and z-scores are reported. Log-transformed response times on correct trials were modeled using linear regression, and t-values are reported. * $p < .05$; ** $p < .01$, *** $p < .001$.

Performance measure (Figure)	Predictor	β	z / t	p	sig.
Overall accuracy (2A)	Age (lin.)	0.054	3.1	0.0017	**
	Age (qua.)	-0.0014	-3.0	0.0024	**
	Sex	0.0074	0.2	0.82	
Response times (2B)	Age (lin.)	-0.17	-8.4	< 0.001	***
	Age (qua.)	-0.004	-7.4	< 0.001	***
	Sex	0.19	5.1	< 0.001	***
Stay after (pot.) switch (2C)	Age (lin.)	0.42	3.8	< 0.001	***
	Age (qua.)	-0.010	-3.5	< 0.001	***
	Sex	0.27	1.3	0.19	
Asymptotic performance (2D)	Age (lin.)	0.19	4.2	< 0.001	***
	Age (qua.)	-0.0048	-4.0	< 0.001	***
	Sex	0.025	0.3	0.77	

166 To determine the age of peak performance, we binned participants into equal-
167 sized groups based on age (Methods; suppl. Fig. 13D-F; Fig. 3C-F). Over-
168 all task performance peaked in 13-15 year-olds (mid-adolescence), and declined
169 steeply for both younger and older participants (Fig. 3C). 13-15 year-olds were
170 also more willing to repeat previous actions after single negative outcomes, es-
171 pecially compared to younger children ("stay" on "(potential) switch trials"; Fig.
172 3E). This suggests that 13-15 year olds were most persistent in the face of negative
173 feedback. 13-15 year-olds also performed best during stable task periods without
174 switches, showing the highest accuracy on asymptotic trials, especially compared
175 to younger participants (Fig. 3F).

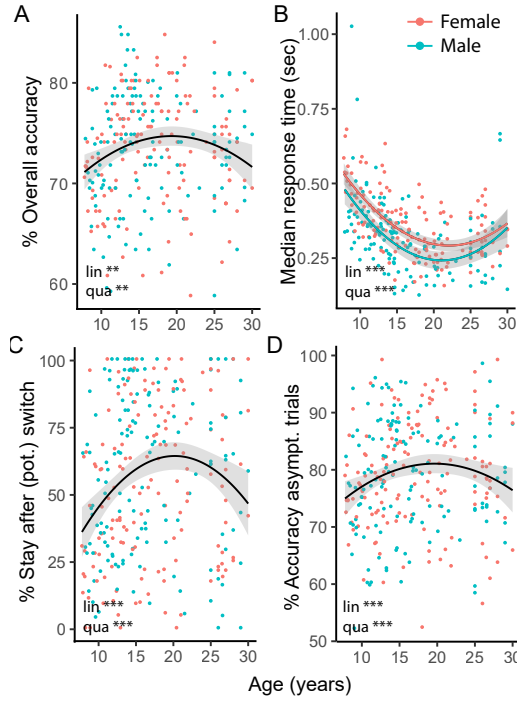


Figure 2: Task performance across age. Each dot shows one participant, color denotes sex. Curves show the fits of linear regression models, with shaded standard errors. "Lin.": significant effect of age on outcome; "qua.": significant effect of squared age on outcome. Stars denote p-values like before. (A) Percentage of correct choices across the entire task (120 trials). (B) Median response times on correct trials. Regression coefficients differed between males and females. (C) Fraction of stay trials after (potential, "pot.") switches (red arrows in Fig. 1C). (D) Accuracy on asymptotic trials (grey bars in Fig. 1C).

176 Furthermore, 13-15 year-olds adapted their choices more optimally to previ-
 177 ous outcomes than younger or older participants. To show this, we used mixed-
 178 effects logistic regression to predict actions on trial t from predictors encoding
 179 positive or negative outcomes on trials $t - i$, for delays $1 \leq i \leq 8$ (Methods). The
 180 effects of positive outcomes were several times larger than the effects of nega-
 181 tive outcomes (suppl. Table 8; Fig. 13B-F), in accordance with task dynamics:
 182 Positive outcomes indicated with certainty that an action was correct, justifying
 183 their strong effect on behavior, whereas negative outcomes were ambivalent as to
 184 whether a switch occurred or not, and should have smaller effects. Crucially, this
 185 pattern differed between participants of different ages, as revealed by interactions
 186 between age and previous outcomes (suppl. Fig. 13B, C, E, and F; suppl. Table

187 8): On trials $t - 1$ and $t - 2$, both positive and negative outcomes interacted with
188 age and squared age (all p 's < 0.014 ; suppl. Table 8), such that the effect of posi-
189 tive outcomes increased with age and then slowly plateaued (suppl. Fig. 13C, F).
190 For negative outcomes, the sign of the interaction was opposite for trials $t - 1$ ver-
191 sus $t - 2$ (all p 's < 0.046 ; suppl. Table 8). This shows that the effect of negative
192 outcomes flipped, being weakest in 13-15 year olds for trial $t - 1$ (Fig. 13F), but
193 strongest for trial $t - 2$. In other words, 13-15 year-olds were best at ignoring sin-
194 gle, ambivalent negative outcomes ($t - 1$), and most likely to integrate long-range,
195 meaningful negative outcomes ($t - 2$), which potentially indicated task switches.
196 To summarize our model-agnostic results, 13-15 year-olds outperformed younger
197 participants 8-13, older adolescents, and adults on a stochastic and volatile task,
198 which was designed to mimic environmental challenges specific to adolescence.
199 We next used computational modeling to investigate what cognitive processes
200 gave rise to 13-15 year old adolescents' superior performance.

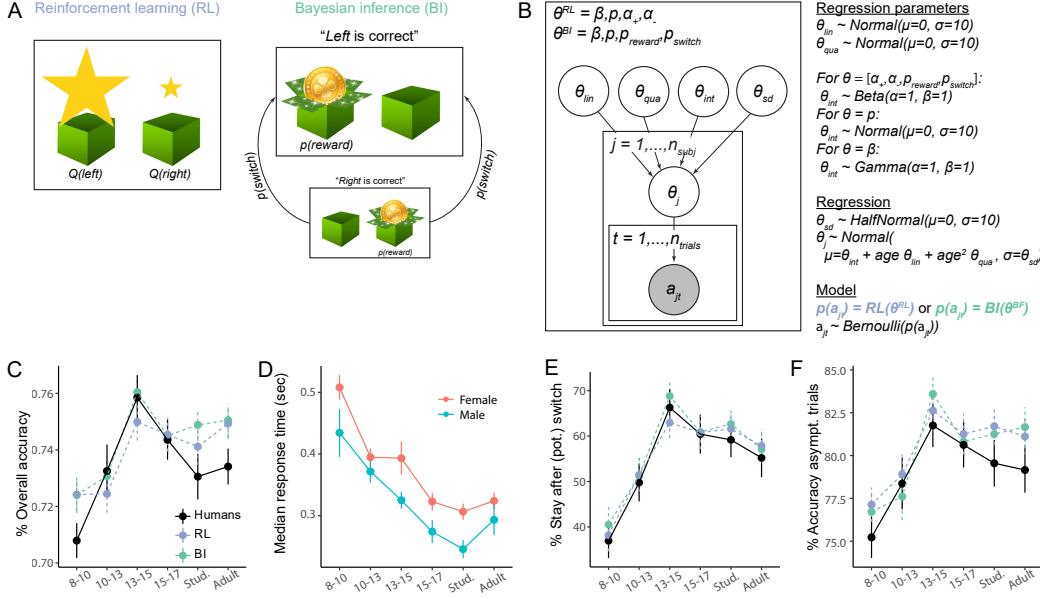


Figure 3: (A) Conceptual depiction of the RL and BI models. In RL (left), actions are selected based on learned values, illustrated by the size of stars ($Q(left)$, $Q(right)$). In BI (right), actions are selected based on a mental model of latent state of the task, which contains task stochasticity ($p(reward)$) and volatility ($p(switch)$). The size of each state illustrates the probability of being in that state. (B) Hierarchical Bayesian model fitting. Box on the left: Models had free parameters θ^{RL} or θ^{BI} . Individual parameters θ_j were based on group-level parameters θ_{sd} , θ_{int} , θ_{lin} , and θ_{qua} (see below). For each model (RL and BI), all parameters were simultaneously fit to the observed (shaded) sequence of actions a_{jt} of all participants j , using MCMC sampling. Right: Priors for group-level parameters were uninformative; the form of the prior differed based on parameter ranges. For each participant j , each parameter θ was sampled according to a linear regression model, based on group-wide standard deviation θ_{sd} , intercept θ_{int} , linear change with age θ_{lin} , and quadratic change with age θ_{qua} . Each model (RL or BI) provided a choice likelihood $p(a_{jt})$ for each participant j on each trial t , based on individual parameters θ_j . Action selection followed a Bernoulli distribution. See Methods for details. (C)-(F) Human behavior on the measures of Fig. 2, binned in age groups. (C), (E), and (F) also show simulated model behavior, verifying that models closely reproduced human behavior and age differences.

3.3. Cognitive Modeling

3.3.1. Models

We fitted two classes of cognitive models to the task, RL and BI. The winning RL model included four parameters: persistence p , inverse decision temperature β , and positive and negative learning rates α_+ and α_- (Methods). Notably, this model updated the values of both the chosen and unchosen action after each out-

come, allowing for counterfactual learning (Boorman et al., 2011; Palminteri et al., 2016). It also allowed learning rates to differ between positive (α_+) and negative outcomes (α_-), an increasingly common idea in cognitive neuroscience (e.g., Frank et al., 2004; van den Bos et al., 2012; for review, see Nussbaum and Hartley, 2019) and AI (Dabney et al., 2020). Parameters p and β controlled the translation from RL values into choices: persistence p increased the probability of repeating choices when $p > 0$, and of alternating choices when $p < 0$; β induced decision noise (increased probability of exploratory choices) when small, and allowed for reward-maximizing choices when large. The winning BI model also had four parameters: choice-parameters p and β as in the RL model, as well as task volatility p_{switch} and reward stochasticity p_{reward} , which characterized participants’ internal model of the task (Fig. 3A; Methods). p_{switch} ranged from stable ($p_{switch} = 0$) to volatile ($p_{switch} > 0$), and p_{reward} ranged from deterministic ($p_{reward} = 1$) to stochastic ($p_{reward} < 1$). The actual task was based on $p_{switch} = 0.05$ and $p_{reward} = 0.75$.

We fitted each model to participant data using hierarchical Bayesian fitting (Fig. 3B; Methods). This approach recovered individual parameters reliably (suppl. Fig. 14), and allowed us to estimate the effects of age on model parameters in a statistically unbiased way (Katahira, 2016; M. D. Lee, 2011; van den Bos et al., 2017). We compared different parameterizations of each model using the WAIC (Watanabe, 2013) to identify a winning RL and a winning BI model (Table 2). The winning RL model had the lowest score overall, revealing best quantitative fit. Nevertheless, both RL and BI models validated equally well, closely reproducing human behavior and age-related differences: Both models showed the performance peak in 13-15 year olds (Fig. 3C), the largest proportion of staying after (potential) switch trials (Fig. 3E), best asymptotic performance on non-switch trials (Fig. 3F), and the most efficient use of previous outcomes to adjust future actions (suppl. Fig. 13 D-F). Other tested models (Table 2) did not capture all qualitative patterns (suppl. Fig. 16, 17). To conclude, despite major differences in their theoretical framework, both RL and BI captured human behavior and age differences. This finding has interesting implications, which we discuss in detail in the Discussion.

Table 2: WAIC model fits and standard errors for all models, based on hierarchical Bayesian fitting. Bold numbers highlight the winning model of each class. For the parameter-free BI model, the Akaike Information Criterion (AIC) was calculated precisely. WAIC differences are relative to next-best model of the same class, and include estimated standard errors of the difference as an indicator of meaningful difference. In the RL model, " α " refers to the classic RL formulation in which $\alpha_+ = \alpha_-$. " α_c " refers to the counterfactual learning rate that guides updates of unchosen actions, with $\alpha_{+c} = \alpha_{-c}$ (Methods).

	Free parameters (count)		(W)AIC	WAIC Difference
BI	–	(0)	31,959	2,668 \pm 0
	β	(1)	29,291 \pm 206	868 \pm 78
	β, p	(2)	28,423 \pm 201	4,769 \pm 132
	β, p, p_{reward}	(3)	23,654 \pm 203	51 \pm 10
	$\beta, p, p_{reward}, p_{switch}$	(4)	23,603 \pm 200	0
RL	α, β	(2)	26,678 \pm 200	438 \pm 44
	α, β, α_c	(3)	26,240 \pm 201	1,429 \pm 78
	$\alpha, \beta, \alpha_c, p$	(4)	24,811 \pm 190	42 \pm 13
	$\alpha_+, \beta, \alpha_{+c}, p, \alpha_-$	(5)	24,769 \pm 213	1,260 \pm 73
	$\alpha_+, \beta, \alpha_{+c}, p, \alpha_-, \alpha_{-c}$	(6)	23,509 \pm 211	17 \pm 10
	$\alpha_+ = \alpha_{+c}, \alpha_- = \alpha_{-c}, \beta, p$	(4)	23,492 \pm 201	0

3.3.2. Age Differences in Model Parameters

All model parameters showed age effects (Fig. 4; suppl. Tables 11 and 12). We tested these effects statistically by modeling age explicitly in a hierarchical Bayesian model (Fig. 3B, suppl. Table 11), and also by assessing age-group differences in the posteriors of an age-less hierarchical Bayesian model (suppl. Table 12; Methods).

Choice-based parameters p and β were almost perfectly correlated between the winning RL and BI models, even though they were fitted independently (Spearman $\rho = 0.94$; Fig. 5B). This suggests that the parameters captured robust, update-independent aspects of decision making. p and β both increased monotonically with age and plateaued in older participants (Fig. 4A, B, E, F). This was reflected in linear and negative quadratic effects of age (suppl. Table 11): Persistence p increased near-linearly from age 8 until 17, and then plateaued around age 18-30 (Fig. 4A, E). This shows that the willingness to repeat previous actions, independent of outcomes, increased from childhood to adulthood, with steady growth during teen years.

Other parameters showed non-monotonic age trajectories. α_- , p_{reward} , and

p_{switch} declined drastically from age 8 to 13-15, but then reversed their trajectory and increased again, reaching a plateau that lasted from 15-30 years (Fig. 4C, G-H). For α_- and p_{reward} , these changes were captured in significant pairwise differences between children (8-10) and 13-15 year-olds, as well as between 13-15 year-olds and adults (25-30; for statistics, see suppl. Table 12; also see Methods). For p_{switch} , age differences were captured in a significant quadratic effect of age (suppl. Table 11). Parameters p_{reward} and p_{switch} , reflecting participants' mental model of the task, were closest to their true values ($p_{reward} = 0.75$; $p_{switch} = 0.05$) in 13-15 year-olds. 8-10 year-old children and adults (18-30) overestimated task volatility (p_{switch}) and underestimated the reward stochasticity (p_{reward}) to a larger degree. Parameter α_- also was lowest in 13-15 year-olds, allowing them to avoid premature switching based on single negative outcomes while allowing for slow integration of outcomes and adaptive switching after multiple negative outcomes. Parameter α_+ showed a unique age trajectory with relatively stable values during childhood and adolescence (8-17), and a sudden increase in adults (18-30; Fig. 4D), captured in a linear effect of age (suppl. Table 11).

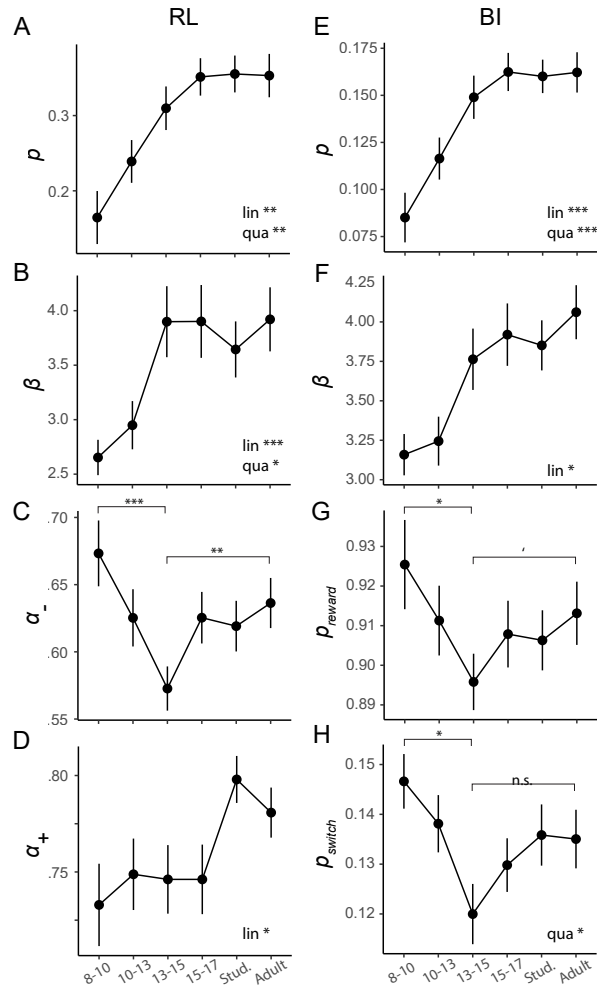


Figure 4: Age-related differences in model parameters for the winning RL (left column) and BI model (right). Stars indicate significant linear ("lin") and quadratic ("qua") effects of age on model parameters, obtained from the hierarchical Bayesian model, and differences between age groups, obtained from the age-less hierarchical Bayesian model (Methods; suppl. Tables 11 and 12). Means and standard errors were calculated based on individual fits from the age-less model, to avoid double-dipping (Methods). (A)-(D) RL model parameters. (E)-(H) BI model parameters. Stars indicate p-values like before.

3.4. Differences between RL and BI

Having obtained two independent sets of parameters for each participant from two computational models based on different cognitive mechanisms, we aimed to

275 clarify how both models were related. We first asked whether each model captured
276 different aspects of behavior, or whether both models captured the same behav-
277 iors and merely differed in form. To test this, we simulated artificial behavior
278 from each model and assessed how well these data were captured by the opposite
279 model. Each model was fitted worse by the opposite model than by itself (Fig.
280 5A), which reveals that each model captured unique aspects of behavior. (This
281 difference was smaller when fitting the RL model, suggesting that it was more
282 versatile and captured more aspects of the BI model than the other way around.)

283 We next asked how closely individual parameters were related between mod-
284 els, assessing pairwise Spearman correlations. As mentioned before, choice pa-
285 rameters p and β were almost perfectly correlated between models (p : $\rho = 0.97$;
286 β : $\rho = 0.94$; Fig. 5B). In addition, parameter p_{reward} (BI) was strongly correlated
287 with α_- (RL), suggesting that negative learning rate (α_-) and beliefs about task
288 stochasticity (p_{reward}) played similar roles in the integration of negative outcomes.
289 Parameter p_{switch} (BI) was strongly negatively correlated with β (RL), suggesting
290 that decision noise (β) in the RL model captured aspects that were explained by
291 beliefs about task volatility (p_{switch}) in the BI model. The only parameter that
292 showed no large correlations with other parameters was α_+ (RL), suggesting a
293 unique role.

294 Lastly, we investigated how much information each model provided about the
295 other, using linear regression to predict each parameter from the parameters and
296 one-way parameter interactions of the other model. Seven out of eight param-
297 eters were predicted almost perfectly (Fig. 5C), showing that the parameters of
298 one model captured almost all variance in the opposite model. In other words,
299 fitting the RL model on participants' data allowed us to nearly perfectly predict
300 participants' BI parameters, without fitting the BI model. Note that α_+ (RL) was
301 again an exception in that its variance was not fully captured by BI parameters.
302 α_+ might thus account for the better fit of the RL model to human (Table 2) and
303 simulated data (Fig. 5A), compared to the BI model.

304 In summary, RL and BI models captured similar aspects of behavior, as shown
305 by large inter-model parameter correlations and amounts of explained variance;
306 nevertheless, both models were not redundant, as evident in the fact that each was
307 unable to perfectly fit the other.

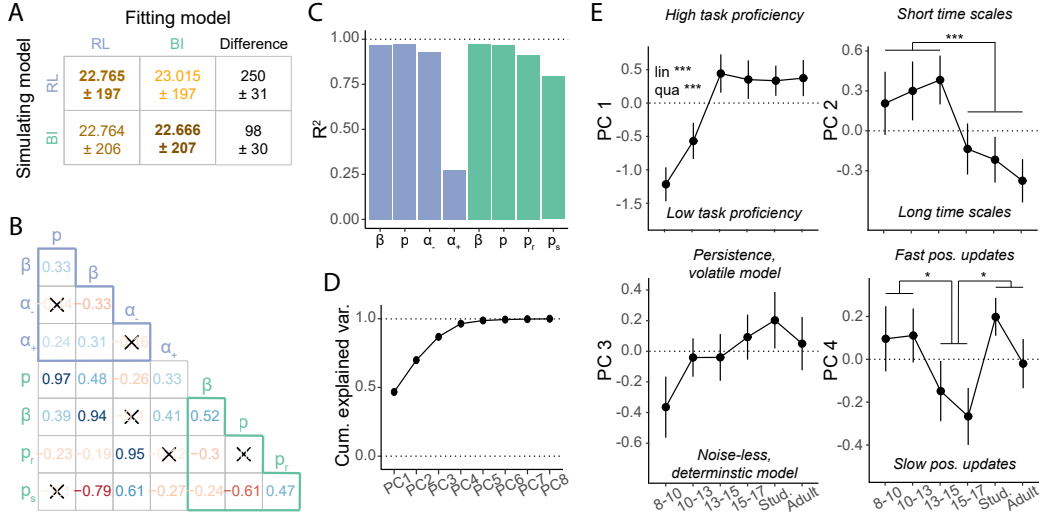


Figure 5: Relating RL and BI models. (A) Model recovery. WAIC scores were worse (larger; lighter colors) when recovering behavior that was simulated from one model (row) using the other model (column), than when using the same model (diagonal), revealing that the models were discriminable. The difference in fit was smaller for BI simulations (bottom row), suggesting that the RL model captured BI behavior better than the other way around (top row). (B) Spearman pairwise correlations between model parameters. Red (blue) hue indicates negative (positive) correlation, saturation indicates correlation strength. Non-significant correlations are crossed out (Bonferroni-corrected at $p = 0.00089$). Light-blue (teal) letters refer to RL (BI) model parameters. Light-blue / teal-colored triangles show correlations within each model, remaining cells show correlations between models. (C) Variance of each parameter explained by parameters and interactions of the other model (R^2), estimated through linear regression. All four BI parameters (green) were predicted almost perfectly by the RL parameters, and all RL parameters except for α_+ (RL) were predicted by the BI parameters. (D)-(E) Results of PCA on model parameters. (D) Cumulative variance explained by all principal components PC1-8. The first four components captured 96.5% of total parameter variance. (E) Age-related differences in PC1-4: PC1 reflected overall task proficiency and showed rapid development between ages 8-13, which were captured by linear ("lin") and quadratic ("qua") effects in a regression model. PC2 captured a step-like transition from shorter to longer updating time scales at age 15, as revealed by PC-based model simulations (Supplements). PC3 showed no significant age effects. PC4 captured the variance in α_+ and differed between adolescents 15-17 and both 8-13 year olds and adults. PC2 and PC4 were analyzed using t-tests. * $p < .05$; ** $p < .01$, *** $p < .001$.

3.5. Combining RL and BI using PCA

We next asked whether both models in conjunction provided a even better explanation of unique adolescent decision making than either model on its own, using PCA to unveil the lower-dimensional structure embedded in the shared 8-

dimensional parameter space. Indeed, only four dimensions were necessary to capture 96.5% of the parameter variance (Fig. 5D), suggesting that individual differences between participants could be explained by variation on just these four, rather than eight different model parameters.

To identify the role of each PC, we assessed parameter loadings and gained additional insight by simulating behavior with parameter sets defined by each PC. This approach is similar to simulating behavior based on different parameter values in order to investigate the effects of parameters (e.g., simulating big versus small values of β produces high versus low task performance). We similarly simulated behavior based on small versus large values of each PC, and compared the resulting behaviors to identify the exact function of each PC (Methods). Principal component 1 (PC1), the dimension capturing the largest proportion of variance, reflected general task proficiency (suppl. Fig. 18A; suppl. Table 13; suppl. Text). Low proficiency was caused by larger-than-average values of α_- (RL), which led to premature switching, and p_{reward} and p_{switch} (BI), reflecting overly deterministic but volatile mental models of the task (suppl. Fig. 18A, left). High proficiency, on the other hand, was caused by larger-than-average values of α_+ (RL), p (RL and BI), and β (RL and BI), which facilitated quick integration of positive outcomes, choice persistence, and low decision noise, respectively (suppl. Fig. 18A, right). PC1 was lowest in the youngest participants (8-10), but increased rapidly until age 13, at which age it reached a stable plateau that lasted throughout adulthood (Fig. 5E, top-left). Age differences in PC1 were characterized by linear and quadratic effects of age (Methods). Taken together, PC1 explained one side of the inverse U-shape in overall task performance (Fig. 2; suppl. Fig. 13; Fig. 3C-F): 13-15 year olds outperformed 8-13 year olds because younger participants had not yet mastered task proficiency.

PC1 did not explain, however, how participants aged 13-15 outperformed older participants. PC2, the dimension that captured the second-most variance after PC1, played this role by capturing the tension between updates with short versus long time scales (suppl. Fig. 18E). Short time scales were driven by larger-than-average values of α_+ and α_- (RL), i.e., rapid updates based on recent outcomes, which led to pronounced win-stay lose-shift behavior (suppl. Fig. 18B, left). In the BI model, short time scales were driven by increasingly volatile (p_{switch}) and especially deterministic (p_{reward}) mental models. Short time scales were complemented by persistence, driven by larger-than-average values of p (RL and BI). Long-term updates were the result of lower-than-average values on all parameters (suppl. Fig. 18B, left), allowing for gradual and slow, but precise value updates, with choice unbiased by persistence. PC2 showed a step function

of age. Whereas younger participants persisted and acted on short times scales, after age 15, participants showed unbiased long-term updates (suppl. Fig. 5E, top-right). Differences between 13-15 year olds and older participants were therefore captured by PC2, suggesting that better performance was the result of shorter updating time scales. This is conditioned on high task proficiency (PC1), which includes slower time scales for negative outcomes, but faster time scales for positive outcomes.

PC3 captured the tension between noise-less choice (larger-than-average β) combined with a deterministic mental model (larger-than-average p_{reward}) on one side, and persistent choice (larger-than-average p) combined with an overly volatile mental model (larger-than-average p_{switch}) on the other (suppl. Fig. 18C). PC3 showed no significant age-related differences (Fig. 5E, bottom-left). PC4 captured the unique variance of α_+ (RL), with a tension between slow and fast updates from positive outcomes (suppl. Fig. 18D). PC4 was lower in 13-17 year olds than both 8-13 year olds and adults (18-30; Fig. 5E, bottom-right), revealing that after accounting for variance in PC1-3, the remaining variance was explained by adolescents' relatively longer updating timescales for positive outcomes. In other words, positive outcomes had weaker immediate, but stronger long-lasting effects in 13-17 year olds, setting them apart from both younger and older participants.

4. Discussion

Across species, the adolescent transition from childhood to adulthood brings great challenges for learning and exploration. From an evolutionary perspective, these challenges may have caused the adolescent brain to evolve behavioral tendencies that promote adaptive learning in rapidly changing, uncertain environments. To test this idea, we examined choice behavior in a stochastic and volatile task adapted from rodent studies (Tai et al., 2012).

Indeed, we found that 13-15 year olds performed better than both younger (8-13) and older participants, including adults (15-30): 13-15 year olds achieved the highest overall accuracy, were most willing to wait out negative feedback (potential switch), and made the best choices during stable periods (asymptotic performance). Overall, 13-15 year olds used negative feedback most optimally to guide future choices, being least affected by proximal, but most sensitive to distal outcomes. This shows an ability to ignore ambivalent information while responding appropriately to meaningful patterns. Indeed, such inverted-U or U-shaped developmental trajectories are not unique in the development of human cognition. Evidence is accumulating that adolescents outperform adults in vari-

ous domains, including probabilistic learning (Davidow et al., 2016) and reversal (van der Schaaf et al., 2011), creativity (Kleibeuker et al., 2013), and social learning (Gopnik et al., 2017). Prowess in flexibility has also been reported in studies of developing rodents (Guskjolen et al., 2017; Johnson and Wilbrecht, 2011; Simon et al., 2013).

One aspect of adolescent behavior, 13-15 year old’s increased willingness to wait out negative feedback, deserves specific attention: It means that in the context of this task, 13-15 year olds were less impulsive than all other age groups. This finding seems inconsistent with past research that often described the mid-adolescent period in terms of increased risk taking and higher risk of negative life outcomes. To explain why, studies of adolescent development have separated impulsivity and sensation seeking. Studies using self reports and experimental tasks showed that impulse control grows through the teen years, while sensation seeking peaks in mid to late adolescence (Albert et al., 2013; Harden and Tucker-Drob, 2011; Romer and Hennessy, 2007; Steinberg et al., 2009). The combination of not-yet-mature impulse control and high sensation seeking in mid-adolescence has been used to explain why this period is associated with higher risks (Harden and Tucker-Drob, 2011; Steinberg, 2013). Our findings do not fit into this narrative of adolescents as risk takers. There are several potential reasons for this discrepancy: (1) Our task may not tap into sensation seeking, a process separate from impulsivity. (2) We use a behavioral task and not self-report methods, which each have different benefits and limitations. (3) Our task may elicit different learning and decision-making strategies than other tasks because it creates a stochastic and volatile environment. Individuals likely do not apply the same learning and decision rules in all contexts. This suggests that the interplay between brain development and the statistics of specific environments may be more important than previously realized (Nussenbaum and Hartley, 2019).

To understand which cognitive and neural processes supported 13-15 year-olds’ superior performance within our specific task, we employed two types of cognitive models, RL and BI. To fit human behavior, the RL model required the ability to learn from counterfactual outcomes (updating values of non-chosen actions), and to apply different learning rates to positive versus negative outcomes (learning parameters α_+ and α_-). It also required persistence, i.e., a tendency to repeat previous actions independent of outcomes, in addition to decision noise, i.e., the ability to explore non-maximizing actions (choice parameters p and β). RL models have been used extensively to shed light on neural mechanisms, and a specialized network of brain regions—including basal ganglia, cortical, and limbic regions—is thought to implement key RL computations (for reviews, see

Frank and Claus, 2006; D. Lee et al., 2012; Niv, 2009; O’Doherty et al., 2015). Fitting RL models to developmental samples is thought to inform our understanding of brain development (e.g., Christakou et al., 2013; Davidow et al., 2016; Javadi et al., 2014; Master et al., 2019; for reviews, see Nussenbaum and Hartley, 2019; van den Bos et al., 2017). Using RL models, we found that choice parameters (β , p) in our study grew monotonically throughout childhood and adolescence, and only matured in late adolescence / early adulthood. This is consistent with previous developmental modeling studies (Nussenbaum and Hartley, 2019), and with a role for late-developing brain circuits in choice behavior (Giedd et al., 1999; Gogtay et al., 2004; Nussenbaum and Hartley, 2019; Sowell et al., 2003; Toga et al., 2006).

While the developmental trajectories of choice parameters have been highly consistent in the developmental modeling literature, the development of learning-rate parameters has been ambivalent and even contradictory (Nussenbaum and Hartley, 2019). One problem might be that many learning studies that are fit with RL models likely involve a variety of different learning processes, which do not only include striatal incremental learning (Yagishita et al., 2014), based on direct and/or indirect pathways (Hauser et al., 2015), but also hippocampal-based episodic memory (Bornstein and Norman, 2017; Wimmer et al., 2014), and frontal-cortical cognitive control (Badre et al., 2010; Collins and Frank, 2012; Daw et al., 2011). Differences in task contexts and task statistics likely elicit different learning strategies, and recruit different neural processes (Nussenbaum and Hartley, 2019). This potentially explains the diversity of previous findings with regard to learning parameters and limits our ability to make inferences about brain development from behavioral modeling studies. It is a future challenge to disentangle the development of multiple systems and context-based responses, for example by studying the same individuals across multiple tasks and computational models.

Another potential reason for the observed discrepancies in learning rates is that between studies, models often differ in the number and type of learning-rate parameters (e.g., positive, negative, factual, counter-factual). For example, the standard learning rate parameter α controls updates to the values of chosen actions after both positive and negative outcomes, whereas a more specialized parameter α_c – controls updates to values of unchosen actions, but only after negative outcomes. Given the likely differences in neural substrate that underlie these different mechanisms, they likely differ in their developmental trajectory. Another reason for the discrepancies in the literature might therefore be study-by-study differences in parameterizations of computational models. In our study, learning rates

462 from negative feedback (α_-) showed a pronounced U pattern with minimum in
 463 13-15 year olds, whereas learning rates from positive feedback (α_+) were stable
 464 throughout childhood and adolescence, then suddenly increased in adults. These
 465 patterns likely reflect the combination of different cognitive and neural processes,
 466 which matured at different times. Indeed, the developmental trajectory of α_- was
 467 almost identical to the stochasticity parameter p_{reward} of the BI model ($\rho = 0.95$),
 468 suggesting that α_- played a role in switching behavior after negative feedback,
 469 rather than learning. As a whole, the RL model might have approximated infer-
 470 ential reasoning rather than performing pure incremental learning, an issue we
 471 discuss in more detail below.

472 To place our RL results in a broader context, we also applied Bayesian Infer-
 473 ence (BI) models to our task data. Using BI models, we found that choice pa-
 474 rameters β and p showed almost identical trajectories as in the RL model, strong
 475 independent support for our hypothesis that these factors increased through the
 476 second decade of life. BI model fits also provided novel results. We found that
 477 BI mental model parameters, task stochasticity p_{reward} and volatility p_{switch} , were
 478 most accurate in 13-15 year-olds. By definition, this means that 13-15 year-olds
 479 possessed the best mental model with respect to actual task statistics, whereas
 480 both younger (8-13) and older participants (15-30) demonstrated less accurate
 481 models. We had also hypothesized that children and adults would expect less
 482 volatility and stochasticity than adolescents, but only stochasticity showed this
 483 pattern. Volatility, on the other hand, appeared to be perceived as larger in younger
 484 (8-13) and older (15-30) participants compared to 13-15 year-olds. Interestingly,
 485 the BI model revealed that participants of all ages deviated markedly from Bayes-
 486 optimal behavior, employing mental models that were both too volatile and too
 487 deterministic (p_{switch} : 8-13 year-olds behaved as if they expected switches every
 488 6.5 trials, 13-15 year-olds every 10 trials, whereas the task switched every 20 trials
 489 on average; p_{reward} : 8-13 year-olds behaved as if they expected rewards for 92%
 490 of correct responses, 13-15 year-olds for 89%, the task rewarded 75%). In sum-
 491 mary, 13-15 year-olds exhibited mental models that were most in line with task
 492 parameters, expecting the most stochasticity and least volatility of all age groups.
 493 This is in accordance with the differentiation between "adaptation" and "settings"
 494 in Nussenbaum and Hartley, 2019, and suggests that 13-15 year-olds showed the
 495 largest ability to adjust and adapt to specific task statistics, rather than reflecting
 496 a particular, developmentally-fixed setting of any specific parameter. A similar
 497 argument about parameter optimality rather than a developmentally-specific pa-
 498 rameter setting was made in Davidow et al., 2016.

499 Fitting two separate model classes led to several benefits in the understand-

ing of the underlying cognitive processes: (1) Both models provided converging (choice parameters) and additive evidence (RL: learning parameters; BI: mental-model parameters). Converging results showed surprisingly strong, direct replication ($\beta_{RL} \leftrightarrow \beta_{BI}$, $p_{RL} \leftrightarrow p_{BI}$), and parallelism between model parameters helped clarify the role of ambiguous parameters ($p_{reward} \rightarrow \alpha_-$). Independent components led to additive insights (e.g., unique parameter α_+), (2) Each model’s conceptual framework and interpretation of the cognitive process became more distinctive in direct comparison with each other. The contrast helped sharpen claims about incremental learning (RL) versus mental-model based inferential reasoning (BI). Whereas the RL model achieved better numerical fit, the BI model provided advantages in terms of interpretability: Our concepts of interest (stochasticity, volatility) were explicitly modeled within the BI framework, potentially allowing insight into how they were processed by each participant. On the technical side, all BI parameters occupied meaningful and interpretable ranges. In the RL model, on the other hand, learning rates showed values substantially larger than 0.5. This made their interpretation difficult, as it is unlikely that high learning rates reflect the type of RL processes implemented in the brain’s RL network, and suggests the model approximated some other dynamic adaptation process. This highlights the fact that numerical model fit (RL) and interpretability (BI) can sometimes be at odds. Future research is necessary to explore this topic in more depth.

A final advantage of fitting both models was the possibility to investigate patterns that go beyond model-specific parameters, using PCA on the shared parameter space. This analysis exposed a different set of factors, which differentiated 13-15 year-olds from younger participants (PC1), from older participants (PC2), or from both (PC4). PC1 reflected overall task proficiency and showed steep improvement until age 13, plateauing thereafter. This suggests that 13-15 year-olds outperformed younger participants because the younger group was too exploratory for the task, was less persistent, and possessed less accurate mental models, leading to weighing negative outcomes too much relative to positive ones. PC2 reflected participants’ updating time scales and showed a step function with transition around 15 years of age. This PC suggests that 13-15 year-olds outperformed older participants because the 15-30 year-olds operated on longer time scales, i.e., were more sensitive to distant outcomes and perceived the task as less stochastic and volatile. PC4 reflected the variance in α_+ that was not captured by previous PCs, and showed an inverse-U trajectory with minimum in 13-17 year-olds. PC4 therefore showed that 13-17 year-olds used the longest time scales when processing positive outcomes, compared to both younger (8-13) and older participants (18-30). Taken together, adolescents aged 13-15 may be at a

”sweet spot” for stochastic and volatile environments because they combine mature levels of task proficiency (PC1) with youthful short updating times scales for all outcomes (PC2), but uniquely long updating time scales for positive outcomes (PC4). This combination would not be optimal in all environments, but in a stochastic and volatile environment, it led to more rewards earned. In this sense, performance in this task supports the idea that the adolescent human brain may pass through stages that have evolved to enhance success in uncertain and volatile environments.

This study shows that age played a crucial role for reward-based decision making and learning in a volatile, stochastic environment. Nevertheless, the question remains which mechanisms underlie these age effects. There is growing evidence that gonadal hormones affect inhibitory neurotransmission, spine pruning, and other variables in the prefrontal cortex of rodents (Delevich et al., 2019; Delevich et al., 2018; Drzewiecki et al., 2016; Juraska and Willing, 2017; Piekarski, Boivin, et al., 2017; Piekarski, Johnson, et al., 2017), suggesting that puberty-related changes in brain chemistry might be the mechanism behind the observed differences. To answer this question, we investigated the trajectories of behavioral performance and model parameters over pubertal development, observing qualitatively similar patterns compared to age (suppl. Fig. 7, 8, 9; suppl. Tables 4, 5; for a discussion of differences, see suppl. Text). Nevertheless, pubertal measures were so highly correlated with age (suppl. Fig. 6) that it was difficult to interpret these findings. We therefore investigated the effects of puberty controlling for age, testing puberty effects separately within each age bin. Puberty effects in this analysis did not reach statistical significance (suppl. Fig. 10, 11, 12). Thus, we were unable to identify a biological mechanism underlying age besides accumulated experience over time.

In conclusion, we used a simple task based on volatility and stochasticity to show that adolescents outperformed adults in a task that represented the kind of learning challenge that may have ecological validity to the transitions and challenges of adolescence. In our community sample, behavior was most optimal at age 13-15. We used two models to examine the underlying cognitive processes. The results suggest that adolescent brains achieved better performance for several reasons: (1) 13-15 year-olds lay on the right spot in a monotonic trajectory between childhood and adulthood (p and β). (2) 13-15 year-olds were outliers in terms of their ability to accurately assess the volatility and stochasticity of their environment, and in terms of their integration of negative outcomes (p_{reward} , p_{switch} , and α_-). (3) 13-15 year-olds combined adult-like (PC1), child-like (PC2), and developmentally unique (PC4) strategies. These data suggest that multiple neu-

576 ral systems underlie developmental changes in brain function, at staggered time
577 scales. Pubertal development and steroid hormones may impact a subset of these
578 processes, yet causality is difficult to determine without manipulation or longitu-
579 dinal designs (Kraemer et al., 2000).

580 For purposes of translation from the lab to the 'real' world, our study indicates
581 that how youth learn and decide changes in a nonlinear fashion as they grow. This
582 underscores the importance of youth-serving programs that are developmentally
583 informed and avoid a one-size-fits-all approach. Finally, these data support a posi-
584 tive view of adolescence and the idea that the adolescent brain exhibits remarkable
585 learning capacities that should be celebrated.

586 **5. Methods**

587 *5.1. Participants*

588 All procedures were approved by the Committee for the Protection of Hu-
589 man Subjects at the University of California, Berkeley. We tested 312 partici-
590 pants: 191 children and adolescents (ages 8-17) and 55 adults (ages 25-30) were
591 recruited from the community and completed a battery of computerized tasks,
592 questionnaires, and saliva samples; 66 university undergraduate students (aged
593 18-50) completed the four tasks as well, but not the questionnaires or saliva sam-
594 ple. Community participants were prescreened for the absence of present or past
595 psychological and neurological disorders; the undergraduate sample indicated the
596 absence of these. Compensation for community participants consisted in 25\$ for
597 the 1-2 hour in-lab portion of the experiment and 25\$ for completing optional
598 take-home saliva samples; undergraduate students received course credit for par-
599 ticipation in the 1-hour study.

600 *Exclusion Criteria.* Out of the 191 participants under 18, 184 completed the stochas-
601 tic switching task; reasons for not completing the task included getting tired, run-
602 ning out of time, and technical issues. Five participants (mean age 10.0 years)
603 were excluded because their mean accuracy was below 58% (chance: 50%), an
604 elbow point in accuracy, which suggests that they did not pay attention to the task.
605 This led to a sample of 179 participants under 18 (male: 96, female: 83). Two par-
606 ticipants from the undergraduate sample were excluded because they were older
607 than 30, leading to a sample aged 18-28; 7 were excluded because they failed
608 to indicate their age. This led to a final sample of 57 undergraduate participants
609 (male: 19, female: 38). All 55 adult community participants (male: 26, female:
610 29) completed the task and were included in the analyses, leading to a sample size

611 of 179 participants below 18, and 291 in total (suppl. Fig. 6). For some analy-
612 ses, we split participants into quantiles based on age. Quantiles were calculated
613 separately within each sex.

614 5.2. *Testing Procedure*

615 After entering the testing room, participants under 18 years and their guardians
616 provided informed assent and permission; participants over 18 provided informed
617 consent. Guardians and participants over 18 filled out a demographic form. Par-
618 ticipants were led into a quiet testing room in view of their guardians, where they
619 used a video game controller to complete four computerized tasks. At the conclu-
620 sion of the tasks, participants between 11 and 18 completed the PDS questionnaire
621 themselves and were measured in height and weight. Participants were then com-
622 pensated with \$25 Amazon gift cards.

623 5.3. *Experimental Design*

624 The task described in this work was the last of the four tasks, a stochastic
625 switching task. The other tasks will be or have been reported elsewhere (Master
626 et al., 2019; Xia et al., 2020). The goal of the stochastic switching task was to
627 collect golden coins, which were hidden in one of two green boxes. On each
628 trial, participants decided which box to open, and task contingencies switched
629 unpredictably throughout the task (Fig. 1B). Before the main task, participants
630 completed a 3-step tutorial: A first prompt explained that one of the two boxes
631 contained a coin (was "magical"), whereas the other one did not. Ten practice
632 trials followed on which one box revealed a coin when selected, whereas the other
633 was empty (deterministic tutorial). The second prompt stated that the magical box
634 would sometimes switch sides. Participants then received eight trials on which the
635 second box contained the coin (but not the first), followed by eight more trials on
636 which the first box contained the coin (but not the second; switching tutorial). The
637 third and last prompt explained that even the magical box did not always contain a
638 coin. This prompt directly led into the main task (stochastic switching), with 120
639 trials.

640 In the main task, the correct box was rewarded in 75% of trials; the incorrect
641 box was never rewarded. After participants reached a performance criterion (see
642 below), it became possible for contingencies to switch (without notice), such that
643 the previously incorrect box was now the correct one. The performance criterion
644 was to collect 7-15 rewards, whereby the specific number was pre-randomized for
645 each block. Any number of non-rewarded trials was allowed in-between rewarded

646 trials. Due to this design, switches only occurred after rewarded trials. For consis-
 647 tency with the rodent version of the task (Tai et al., 2012), the first correct choice
 648 after a switch was always rewarded (not just in %75).

649 5.4. Behavioral Analyses

650 We assessed the effects of age on behavioral outcomes (Fig. 2), using (logistic)
 651 mixed-effects regression models with the package lme4 (Bates et al., 2015) in R
 652 (RCoreTeam, 2016). All models included the following set of regressors to predict
 653 outcomes of interest (e.g., overall accuracy, response times): Age, to assess the
 654 linear effect of age on the outcome; squared age, to assess the quadratic effect
 655 of age; and sex; furthermore all models specified random effects of participants,
 656 allowing participants’ intercepts and slopes to vary independently. When models
 657 included additional predictors, this is noted in the main text.

658 We assessed the effects of previous outcomes on participants’ choices (suppl.
 659 Fig. 13B, C, E, F) using a logistic mixed-effects regression model, which pre-
 660 dicted actions (left, right) from previous outcomes (details below), while testing
 661 for effects of and interactions with sex, z-scored age, and z-scored quadratic age,
 662 specifying participants as mixed effects. We included one predictor for positive
 663 and one for negative outcomes at each delay i with respect to the predicted action
 664 (e.g., $i = 1$ trial ago). Outcome predictors were coded -1 for left and +1 for right
 665 choices, and 0 otherwise. Including predictors of trials $1 \leq i \leq 8$ provided the best
 666 model fit (suppl. Table 8). To visualize the results of this grand regression model
 667 (including all participants), we ran a separate model for each participant with the
 668 same structure, and show individual fits in suppl. Fig. 13B, C, E, F.

669 5.5. Computational Models

670 5.5.1. Reinforcement Learning (RL) Models

In RL, decisions are made based on action values, which are continuously
 updated based on outcomes (Sutton and Barto, 2017). A simple RL model has two
 parameters, learning rate α and decision temperature β . On each trial t , the value
 $Q_t(a)$ of action a is updated based on the observed outcome $o_t \in [0, 1]$ (reward, no
 reward), in the following way:

$$Q_{t+1}(a) = Q_t(a) + \alpha(o_t - Q_t(a))$$

671 I.e., previous action values are updated in proportion to the difference between
 672 the estimated value and the actual reward, scaled by the learning rate α . The
 673 difference itself, $o_t - Q_t(a)$, is called "reward prediction error".

Over time, action values approximate the true underlying reward probabilities. Decisions are based on these values by calculating action probabilities using a softmax transform:

$$p_t(a) = \frac{\exp(\beta Q_t(a))}{\exp(\beta Q_t(a)) + \exp(\beta Q_t(a_{ns}))}$$

674 Here, a is the selected, and a_{ns} the non-selected action.

675 The best-fit 4-parameter RL model was based on this 2-parameter model, with
 676 additional parameters learning rate for negative outcomes α_- , persistence p , as
 677 well as counterfactual reasoning (see below). Adding α_- allowed for separate
 678 updates of rewarded ($o_t = 1$) and non-rewarded ($o_t = 0$) trials: $Q_t(a) = Q_t(a) +$
 679 $\alpha_+(o_t - Q_t(a))$ iff $o_t = 1$, and $Q_t(a) = Q_t(a) + \alpha_-(o_t - Q_t(a))$ iff $o_t = 0$, with
 680 independent α_- and α_+ . Choice persistence or "stickiness" p changed the value
 681 of the previously-selected action a_t on the subsequent trial, biasing toward staying
 682 ($p > 0$) or switching ($p < 0$): $Q(a_t) = Q(a_t) + p$ iff $a_t = a_{t-1}$.

683 Counterfactual reasoning was implemented through updates to the values of
 684 non-selected actions, using counterfactual outcomes $1 - o_t$: $Q_{t+1}(a_{ns}) = Q_t(a_{ns}) +$
 685 $\alpha_+((1 - o_t) - Q_t(a_{ns}))$ iff $o = 1$, and $Q_{t+1}(a_{ns}) = Q_t(a_{ns}) + \alpha_-((1 - o_t) - Q_t(a_{ns}))$
 686 iff $o = 0$. Initially, we used separate parameters α_{+c} and α_{-n} for counterfactual
 687 updates, which were independent from α_+ and α_- for factual updates. Never-
 688 theless, collapsing $\alpha_+ = \alpha_{+c}$ and $\alpha_- = \alpha_{-n}$ improved model fit (Table 2). This
 689 shows that outcomes triggered equal-sized updates to chosen and unchosen ac-
 690 tions. Explained differently, the final model based decisions on a single value
 691 estimate—the value difference between the two available actions—, rather than
 692 on an independent value estimates for each. Chosen and unchosen actions were
 693 updated to the same degree and in opposite directions on each trial.

694 Action values were initialized at 0.5 for all models, reflecting equal initial
 695 values for the two actions.

696 5.5.2. Bayesian Inference (BI) Models

The BI model assumes that participants know that the task has two latent states: "Left action is correct" ($a_{left} = cor$) and "Right action is correct" ($a_{right} = cor$), where *cor* stands for correct (*inc*: incorrect). Participants assume that on each trial, the latent state switches with probability p_{switch} , and that in each state, the probability of receiving a reward for the correct action is p_{reward} (Fig. 3A). On each trial, participants select an action in two phases, using the Bayesian Filter algorithm (Sarkka, 2013): (1) In the *estimation phase*, participants infer the hidden

state of the previous trial $t - 1$, based on the outcome o_{t-1} they received for their action a_{t-1} , using Bayes rule:

$$p(a_{t-1} = cor | o_{t-1}) = \frac{p(o_{t-1} | a_{t-1} = cor) p(a_{t-1} = cor)}{p(o_{t-1} | a_{t-1} = cor) p(a_{t-1} = cor) + p(o_{t-1} | a_{t-1} = inc) p(a_{t-1} = inc)}$$

$p(a_{t-1} = cor)$ is the prior probability that a_{t-1} was correct (on the first trial, $p(a = cor) = 0.5$ for both actions), and $p(o_{t-1} | a_{t-1})$ is the likelihood of the observed outcome o_{t-1} given action a_{t-1} . According to the mental model, likelihoods are (dropping underscripts for clarity): $p(o = 1 | a = cor) = p_{reward}$, $p(o = 0 | a = cor) = 1 - p_{reward}$, $p(o = 1 | a = inc) = \epsilon$, and $p(o = 0 | a = inc) = 1 - \epsilon$, where ϵ is the probability of receiving a reward for an incorrect action, which was 0 in reality, but we set $\epsilon = 0.0001$ to avoid model degeneracy. (2) In the *prediction phase*, participants integrate the possibility of state switches by propagating the inferred knowledge about the hidden state at $t - 1$ forward to trial t :

$$p(a_t = cor) = (1 - p_{switch}) p(a_{t-1} = cor) + p_{switch} p(a_{t-1} = inc)$$

697 We first assessed a parameter-free version of the BI model, truthfully setting
698 $p_{reward} = 0.75$, and $p_{switch} = 0.05$. Lacking free parameters, this model
699 was unable to capture individual differences and led to poor qualitative (suppl.
700 Fig. 17A) and quantitative model fit (Table 2). The best-fit BI model had four
701 free parameters: p_{reward} and p_{switch} , as well as the choice parameters β and p ,
702 like the winning RL model. β and p were introduced by applying a softmax
703 to $p(a_t = cor)$ to calculate $p(a_t)$, the probability of selecting action a on trial
704 t : $p(a_t) = \frac{1}{(1 + \exp(\beta(0.5 - p - p(a_t = cor))))}$. When both actions had the same probability
705 $p(a)$ and persistence $p > 0$, then staying was more likely; when $p < 0$, then
706 switching was more likely.

707 5.5.3. Model Fitting and Comparison

We fitted parameters using hierarchical Bayesian methods (M. D. Lee, 2011; Fig. 3B), and found that the obtained results clearly superseded those of classical maximum-likelihood fitting in terms of parameter recovery (suppl. Fig. 14). Hierarchical Bayesian model fitting estimates the parameters of an entire population *data* jointly, using Bayes formula:

$$p(\theta | data) \propto p(data | \theta) p(\theta)$$

Individual parameters are embedded in a hierarchical structure, which helps resolve uncertainty at the individual level. Because we were interested in age-related

differences in model parameters, we used a hierarchical structure in which parameters $\theta_j^{RL} = [p, \beta, \alpha_-, \alpha_+]$ or $\theta_j^{BI} = [p, \beta, p_{switch}, p_{reward}]$ of participant j were embedded in linear regressions:

$$\theta_j \sim \text{Normal}(\mu = \theta_{int} + \text{age} \theta_{lin} + \text{age}^2 \theta_{qua}, \sigma = \theta_{sd})$$

Each parameter θ was characterized by group-level intercept θ_{int} , slope θ_{lin} , and quadratic change with age θ_{qua} . Individual parameters θ_j were drawn from a normal distribution with standard deviation θ_{sd} around this regression line (Fig. 3B).

Because posteriors $p(\theta|data)$ were analytically intractable, we approximated them using Markov-Chain Monte Carlo sampling (no-U-Turn sampler), using the PyMC3 package in python (Salvatier et al., 2016). We ran 2 chains per model with 6,000 samples per chain, discarding the first 1,000 as burn-in. All models converged with small MC errors, sufficient effective sample sizes, and \hat{R} close to 1 (suppl. Table 10). Point estimates for individual parameters θ_j were calculated as the mean over all posterior samples. For model comparison, we used the Watanabe-Akaike information criterion (WAIC), which estimates the expected out-of-sample prediction error using a bias-corrected adjustment of within-sample error (Watanabe, 2013).

To statistically test the hypothesis that parameter θ differed between age groups, we fitted a separate hierarchical Bayesian model, which did not have access to participants' age, called the "age-free" model. Instead of lying on an age-based regression line, all individual parameters were drawn from the same group-wide Normal distribution with mean θ_{mean} and standard deviation θ_{sd} . To test for differences between groups $a1$ and $a2$ without the danger of double-dipping, we assessed $\theta_{a1} < \theta_{a2}$ in each posterior sample of this model, and then calculated $p(\theta_{a1} < \theta_{a2})$ across all samples. The age-less model was also used to visualize individual parameters in suppl. Figures 15, 17, and 16, and to calculate group means in Fig. 4. Using the age-less model avoided double-dipping on age effects, which would occur if we plotted parameters across age that were fitted in an age-dependent model.

5.5.4. Integrating RL and BI Models

Model Recovery between RL and BI (Fig. 5A). We simulated one dataset per participant from each model, using parameters fitted by the age-free model. We then fitted the simulated data with both models using age-free hierarchical Bayesian fitting. We finally calculated WAIC scores and standard errors using PyMC3 (Salvatier et al., 2016).

740 *Correlations between Model Parameters (Fig. 5B).* We used Spearman correlation, the nonparametric version of the Pearson product-moment correlation, because parameters followed different, not necessarily normal, distributions. Results were similar when using Pearson correlation. p-values were corrected for multiple comparisons using the Bonferroni method.

745 *Predicting Parameters from Parameters of the Other Model (Fig. 5C).* We ran eight different regression models, predicting each parameter from the four parameters of the opposite models as well as their one-way interactions, using linear regression in R (RCoreTeam, 2016). Fig. 5C shows the explained variance (R^2) of each model.

750 5.5.5. Principal Component Analysis (PCA)

751 To extract components that covary across parameters, we ran PCA on the fitted parameters data (8 parameters per participant). PCA can be understood as a method that rotates the coordinate system to align the first axis with the dimension of largest variation in the dataset (first principle component; PC), the second axis with the dimension of second-largest variance (second PC), while being orthogonal to the first, and so on. In this way, all resulting PCs are orthogonal to each other, and explain subsequently less variance. We conducted a PCA after centering and scaling (z-scoring) the data, using the statistical programming language R (RCoreTeam, 2016).

760 *Age Differences in PCs (Fig. 5E).* For each PC, we ran similar regression models as for our behavioral measures of performance, predicting participants' PCs from age (linear), age (quadratic), and sex. When significant, effects were noted in Fig. 5E. For PC2 and PC4, we also conducted post-hoc t-tests, correcting for multiple comparison using the Bonferroni method (Table 3).

Table 3: Results of t-tests on PC2 and PC4. df: Welch-adjusted degrees of freedom.

Comparison	<i>t</i>	df	<i>p</i>	Sig.
PC2 (8-15 vs. 15-30)	3.44	266.2	< 0.001	***
PC4 (8-13 vs. 13-17)	2.28	176.8	0.047	*
PC4 (13-17 vs. 18-30)	2.49	176.6	0.028	*

6. Acknowledgments

Numerous people contributed to this research: Lance Kriegsfeld, Celia Ford, Jennifer Pfeifer, Megan Johnson, Gautam Agarwal, Liyu Xia, Rachel Arsenault, Josephine Christon, Shoshana Edelman, Lucy Eletel, Neta Gotlieb, Haley Keglovits, Julie Liu, Justin Morillo, Nithya Rajakumar, Nick Spence, Tanya Smith, Benjamin Tang, Talia Welte, Lucy Whitmore, and Amy Zou. We are also grateful to our participants and their families. The work was funded by National Science Foundation SL-CN grant 1640885 to RD, AGEK, and LW.

References

- Albert, D., Chein, J., & Steinberg, L. (2013). The Teenage Brain: Peer Influences on Adolescent Decision Making. *Current Directions in Psychological Science*, 22(2), 114–120. <https://doi.org/10.1177/0963721412471347>
- Badre, D., Kayser, A. S., & D’Esposito, M. (2010). Frontal Cortex and the Discovery of Abstract Action Rules. *Neuron*, 66(2), 315–326. <https://doi.org/10.1016/j.neuron.2010.03.025>
- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting Linear Mixed-Effects Models Using lme4. *Journal of Statistical Software*, 67(1), 1–48. <https://doi.org/10.18637/jss.v067.i01>
- Boorman, E. D., Behrens, T. E., & Rushworth, M. F. (2011). Counterfactual Choice and Learning in a Neural Network Centered on Human Lateral Frontopolar Cortex (M. L. Platt, Ed.). *PLoS Biology*, 9(6), e1001093. <https://doi.org/10.1371/journal.pbio.1001093>
- Bornstein, A. M., & Norman, K. A. (2017). Reinstated episodic context guides sampling-based decisions for reward. *Nature Neuroscience*, 20(7), 997–1003. <https://doi.org/10.1038/nn.4573>
- Christakou, A., Gershman, S. J., Niv, Y., Simmons, A., Brammer, M., & Rubia, K. (2013). Neural and psychological maturation of decision-making in adolescence and young adulthood. *Journal of Cognitive Neuroscience*, 25(11), 1807–1823. https://doi.org/10.1162/jocn_a_00447
- Collins, A. G., & Frank, M. J. (2012). How much of reinforcement learning is working memory, not reinforcement learning? A behavioral, computational, and neurogenetic analysis: Working memory in reinforcement learning. *European Journal of Neuroscience*, 35(7), 1024–1035. <https://doi.org/10.1111/j.1460-9568.2011.07980.x>

- 799 Dabney, W., Kurth-Nelson, Z., Uchida, N., Starkweather, C. K., Hassabis, D.,
800 Munos, R., & Botvinick, M. (2020). A distributional code for value in
801 dopamine-based reinforcement learning. *Nature*, 577(7792), 671–675. <https://doi.org/10.1038/s41586-019-1924-6>
802
- 803 Dahl, R. E., Allen, N. B., Wilbrecht, L., & Suleiman, A. B. (2018). Importance
804 of investing in adolescence from a developmental science perspective. *Nature*, 554(7693), 441–450. <https://doi.org/10.1038/nature25770>
805
- 806 Davidow, J. Y., Foerde, K., Galvan, A., & Shohamy, D. (2016). An Upside to
807 Reward Sensitivity: The Hippocampus Supports Enhanced Reinforcement
808 Learning in Adolescence. *Neuron*, 92(1), 93–99. <https://doi.org/10.1016/j.neuron.2016.08.031>
809
- 810 Daw, N. D., Gershman, S. J., Seymour, B., Dayan, P., & Dolan, R. J. (2011).
811 Model-Based Influences on Humans’ Choices and Striatal Prediction Errors. *Neuron*, 69(6), 1204–1215. <https://doi.org/10.1016/j.neuron.2011.02.027>
812
813
- 814 Delevich, K., Piekarski, D., & Wilbrecht, L. (2019). Neuroscience: Sex Hormones
815 at Work in the Neocortex. *Current Biology*, 29(4), R122–R125. <https://doi.org/10.1016/j.cub.2019.01.013>
816
- 817 Delevich, K., Thomas, A. W., & Wilbrecht, L. (2018). Adolescence and “Late
818 Blooming” Synapses of the Prefrontal Cortex. *Cold Spring Harbor Symposium on Quantitative Biology*, 83, 37–43. <https://doi.org/10.1101/sqb.2018.83.037507>
819
820
- 821 DePasque, S., & Galván, A. (2017). Frontostriatal development and probabilistic
822 reinforcement learning during adolescence. *Neurobiology of Learning and Memory*, 143, 1–7. <https://doi.org/10.1016/j.nlm.2017.04.009>
823
- 824 Drzewiecki, C. M., Willing, J., & Juraska, J. M. (2016). Synaptic number changes
825 in the medial prefrontal cortex across adolescence in male and female rats:
826 A role for pubertal onset. *Synapse (New York, N.Y.)*, 70(9), 361–368. <https://doi.org/10.1002/syn.21909>
827
- 828 Frank, M. J., & Claus, E. D. (2006). Anatomy of a decision: Striato-orbitofrontal
829 interactions in reinforcement learning, decision making, and reversal. *Psychological Review*, 113(2), 300–326. <https://doi.org/10.1037/0033-295X.113.2.300>
830
831
- 832 Frank, M. J., Seeberger, L. C., & O’Reilly, R. C. (2004). By Carrot or by Stick:
833 Cognitive Reinforcement Learning in Parkinsonism. *Science*, 306(5703),
834 1940–1943. <https://doi.org/10.1126/science.1102941>

- 835 Frankenhuys, W. E., & Walasek, N. (2020). Modeling the evolution of sensitive
836 periods. *Developmental Cognitive Neuroscience*, 41, 100715. <https://doi.org/10.1016/j.dcn.2019.100715>
837
- 838 Giedd, J. N., Blumenthal, J., Jeffries, N. O., Castellanos, F. X., Liu, H., Zijdenbos,
839 A., Paus, T., Evans, A. C., & Rapoport, J. L. (1999). Brain development
840 during childhood and adolescence: A longitudinal MRI study. *Nature Neu-*
841 *roscience*, 2(10), 861–863. <https://doi.org/10.1038/13158>
- 842 Gogtay, N., Giedd, J. N., Lusk, L., Hayashi, K. M., Greenstein, D., Vaituzis, A. C.,
843 Nugent, T. F., Herman, D. H., Clasen, L. S., Toga, A. W., Rapoport, J. L.,
844 & Thompson, P. M. (2004). Dynamic mapping of human cortical devel-
845 opment during childhood through early adulthood. *Proceedings of the Na-*
846 *tional Academy of Sciences*, 101(21), 8174–8179. <https://doi.org/10.1073/pnas.0402680101>
847
- 848 Gopnik, A., O’Grady, S., Lucas, C. G., Griffiths, T. L., Wente, A., Bridgers, S.,
849 Aboody, R., Fung, H., & Dahl, R. E. (2017). Changes in cognitive flexi-
850 bility and hypothesis search across human life history from childhood to
851 adolescence to adulthood. *Proceedings of the National Academy of Sci-*
852 *ences*, 114(30), 7892–7899. <https://doi.org/10.1073/pnas.1700811114>
- 853 Guskjolen, A., Josselyn, S. A., & Frankland, P. W. (2017). Age-dependent changes
854 in spatial memory retention and flexibility in mice. *Neurobiology of Learn-*
855 *ing and Memory*, 143, 59–66. <https://doi.org/10.1016/j.nlm.2016.12.006>
- 856 Harden, K. P., & Tucker-Drob, E. M. (2011). Individual differences in the de-
857 velopment of sensation seeking and impulsivity during adolescence: Fur-
858 ther evidence for a dual systems model. *Developmental Psychology*, 47(3),
859 739–746. <https://doi.org/10.1037/a0023279>
- 860 Hauser, T. U., Iannaccone, R., Walitza, S., Brandeis, D., & Brem, S. (2015).
861 Cognitive flexibility in adolescence: Neural and behavioral mechanisms
862 of reward prediction error processing in adaptive decision making during
863 development. *NeuroImage*, 104, 347–354. <https://doi.org/10.1016/j.neuroimage.2014.09.018>
864
- 865 Javadi, A. H., Schmidt, D. H. K., & Smolka, M. N. (2014). Adolescents adapt
866 more slowly than adults to varying reward contingencies. *Journal of Cog-*
867 *nitive Neuroscience*, 26(12), 2670–2681. [https://doi.org/10.1162/jocn.a.](https://doi.org/10.1162/jocn.a.00677)
868 [00677](https://doi.org/10.1162/jocn.a.00677)
- 869 Johnson, C., & Wilbrecht, L. (2011). Juvenile mice show greater flexibility in mul-
870 tiple choice reversal learning than adults. *Developmental Cognitive Neu-*
871 *roscience*, 1(4), 540–551. <https://doi.org/10.1016/j.dcn.2011.05.008>

- 872 Juraska, J. M., & Willing, J. (2017). Pubertal onset as a critical transition for
873 neural development and cognition. *Brain Research*, 1654(Pt B), 87–94.
874 <https://doi.org/10.1016/j.brainres.2016.04.012>
- 875 Katahira, K. (2016). How hierarchical models improve point estimates of model
876 parameters at the individual level. *Journal of Mathematical Psychology*,
877 73, 37–58. <https://doi.org/10.1016/j.jmp.2016.03.007>
- 878 Kleibeuker, S. W., Dreu, C. K. W. D., & Crone, E. A. (2013). The development of
879 creative cognition across adolescence: Distinct trajectories for insight and
880 divergent thinking. *Developmental Science*, 16(1), 2–12. <https://doi.org/10.1111/j.1467-7687.2012.01176.x>
- 881 Kraemer, H. C., Yesavage, J. A., Taylor, J. L., & Kupfer, D. (2000). How can we
882 learn about developmental processes from cross-sectional studies, or can
883 we? *The American Journal of Psychiatry*, 157(2), 163–171. <https://doi.org/10.1176/appi.ajp.157.2.163>
- 884 Laube, C., Lorenz, R., & van den Bos, W. (2020). Pubertal testosterone corre-
885 lates with adolescent impatience and dorsal striatal activity. *Developmen-*
886 *tal Cognitive Neuroscience*, 42, 100749. <https://doi.org/10.1016/j.dcn.2019.100749>
- 887
888
889
- 890 Lee, D., Seo, H., & Jung, M. W. (2012). Neural Basis of Reinforcement Learning
891 and Decision Making. *Annual review of neuroscience*, 35, 287–308. <https://doi.org/10.1146/annurev-neuro-062111-150512>
- 892
893 Lee, M. D. (2011). How cognitive modeling can benefit from hierarchical Bayesian
894 models. *Journal of Mathematical Psychology*, 55(1), 1–7. <https://doi.org/10.1016/j.jmp.2010.08.013>
- 895
896 Lourenco, F., & Casey, B. (2013). Adjusting behavior to changing environmen-
897 tal demands with development. *Neuroscience & Biobehavioral Reviews*,
898 37(9), 2233–2242. <https://doi.org/10.1016/j.neubiorev.2013.03.003>
- 899 Master, S. L., Eckstein, M. K., Gotlieb, N., Dahl, R., Wilbrecht, L., & Collins,
900 A. G. E. (2019). Distangling the systems contributing to changes in
901 learning during adolescence. *bioRxiv*, 622860. <https://doi.org/10.1101/622860>
- 902
903 Natterson-Horowitz, D. B., & Bowers, K. (2019). *Wildhood: The Astounding Con-*
904 *nections between Human and Animal Adolescents*. New York, Scribner.
- 905 Niv, Y. (2009). Reinforcement learning in the brain. *Journal of Mathematical Psy-*
906 *chology*, 53(3), 139–154.
- 907 Nussenbaum, K., & Hartley, C. A. (2019). Reinforcement learning across devel-
908 opment: What insights can we draw from a decade of research? *Develop-*

909 *mental Cognitive Neuroscience*, 40, 100733. <https://doi.org/10.1016/j.dcn.2019.100733>

910

911 O'Doherty, J. P., Lee, S. W., & McNamee, D. (2015). The structure of reinforcement-

912 learning mechanisms in the human brain. *Current Opinion in Behavioral*

913 *Sciences*, 1, 94–100. <https://doi.org/10.1016/j.cobeha.2014.10.004>

914 Palminteri, S., Kilford, E. J., Coricelli, G., & Blakemore, S.-J. (2016). The Com-

915 putational Development of Reinforcement Learning during Adolescence.

916 *PLoS Computational Biology*, 12(6). <https://doi.org/10.1371/journal.pcbi.1004953>

917

918 Palminteri, S., Wyart, V., & Koechlin, E. (2017). The Importance of Falsifica-

919 tion in Computational Cognitive Modeling. *Trends in Cognitive Sciences*,

920 21(6), 425–433. <https://doi.org/10.1016/j.tics.2017.03.011>

921 Perfors, A., Tenenbaum, J. B., Griffiths, T. L., & Xu, F. (2011). A tutorial intro-

922 duction to Bayesian models of cognitive development. *Cognition*, 120(3),

923 302–321. <https://doi.org/10.1016/j.cognition.2010.11.015>

924 Petersen, A. C., Crockett, L., Richards, M., & Boxer, A. (1988). A self-report

925 measure of pubertal status: Reliability, validity, and initial norms. *Journal*

926 *of Youth and Adolescence*, 17(2), 117–133. <https://doi.org/10.1007/BF01537962>

927

928 Piekarski, D. J., Boivin, J. R., & Wilbrecht, L. (2017). Ovarian Hormones Orga-

929 nize the Maturation of Inhibitory Neurotransmission in the Frontal Cortex

930 at Puberty Onset in Female Mice. *Current biology: CB*, 27(12), 1735–

931 1745.e3. <https://doi.org/10.1016/j.cub.2017.05.027>

932 Piekarski, D. J., Johnson, C. M., Boivin, J. R., Thomas, A. W., Lin, W. C., Dele-

933 vich, K., M Galarce, E., & Wilbrecht, L. (2017). Does puberty mark a

934 transition in sensitive periods for plasticity in the associative neocortex?

935 *Brain Research*, 1654(Pt B), 123–144. <https://doi.org/10.1016/j.brainres.2016.08.042>

936

937 RCoreTeam. (2016). *R: A Language and Environment for Statistical Computing*.

938 Vienna, Austria, R Foundation for Statistical Computing.

939 Romer, D., & Hennessy, M. (2007). A Biosocial-Affect Model of Adolescent Sen-

940 sation Seeking: The Role of Affect Evaluation and Peer-Group Influence

941 in Adolescent Drug Use. *Prevention Science*, 8(2), 89. <https://doi.org/10.1007/s11121-007-0064-7>

942

943 Salvatier, J., Wiecki, T. V., & Fonnesbeck, C. (2016). Probabilistic programming

944 in Python using PyMC3. *PeerJ Computer Science*, 2, e55. <https://doi.org/10.7717/peerj-cs.55>

945

- 946 Sarkka, S. (2013). *Bayesian Filtering and Smoothing*. Cambridge, Cambridge
947 University Press. <https://doi.org/10.1017/CBO9781139344203>
- 948 Sercombe, H. (2014). Risk, adaptation and the functional teenage brain. *Brain and*
949 *Cognition*, 89, 61–69. <https://doi.org/10.1016/j.bandc.2014.01.001>
- 950 Simon, N. W., Gregory, T. A., Wood, J., & Moghaddam, B. (2013). Differences in
951 response initiation and behavioral flexibility between adolescent and adult
952 rats. *Behavioral Neuroscience*, 127(1), 23–32. <https://doi.org/10.1037/a0031328>
- 953
- 954 Sowell, E. R., Peterson, B. S., Thompson, P. M., Welcome, S. E., Henkenius,
955 A. L., & Toga, A. W. (2003). Mapping cortical change across the human
956 life span. *Nature Neuroscience*, 6(3), 309–315. <https://doi.org/10.1038/nrn1008>
- 957
- 958 Steinberg, L. (2005). Cognitive and affective development in adolescence. *Trends*
959 *in Cognitive Sciences*, 9(2), 69–74. <https://doi.org/10.1016/j.tics.2004.12.005>
- 960
- 961 Steinberg, L. (2013). The influence of neuroscience on US Supreme Court de-
962 cisions about adolescents' criminal culpability. *Nature Reviews Neuro-*
963 *science*, 14(7), 513–518. <https://doi.org/10.1038/nrn3509>
- 964 Steinberg, L., Graham, S., O'Brien, L., Woolard, J., Cauffman, E., & Banich,
965 M. (2009). Age Differences in Future Orientation and Delay Discounting.
966 *Child Development*, 80(1), 28–44. <https://doi.org/10.1111/j.1467-8624.2008.01244.x>
- 967
- 968 Sutton, R. S., & Barto, A. G. (2017). *Reinforcement Learning: An Introduction*
969 (2nd ed.). Cambridge, MA; London, England, MIT Press.
- 970 Tai, L.-H., Lee, A. M., Benavidez, N., Bonci, A., & Wilbrecht, L. (2012). Tran-
971 sient stimulation of distinct subpopulations of striatal neurons mimics changes
972 in action value. *Nature Neuroscience*, 15(9), 1281–1289. <https://doi.org/10.1038/nrn.3188>
- 973
- 974 Toga, A. W., Thompson, P. M., & Sowell, E. R. (2006). Mapping brain maturation.
975 *Trends in neurosciences*, 29(3), 148–159. <https://doi.org/10.1016/j.tins.2006.01.007>
- 976
- 977 van den Bos, W., Bruckner, R., Nassar, M. R., Mata, R., & Eppinger, B. (2017).
978 Computational neuroscience across the lifespan: Promises and pitfalls. *De-*
979 *velopmental Cognitive Neuroscience*. <https://doi.org/10.1016/j.dcn.2017.09.008>
- 980
- 981 van den Bos, W., Cohen, M. X., Kahnt, T., & Crone, E. A. (2012). Striatum–Medial
982 Prefrontal Cortex Connectivity Predicts Developmental Changes in Rein-

- 983 forcement Learning. *Cerebral Cortex*, 22(6), 1247–1255. [https://doi.org/](https://doi.org/10.1093/cercor/bhr198)
984 10.1093/cercor/bhr198
- 985 van der Schaaf, M. E., Warmerdam, E., Crone, E. A., & Cools, R. (2011). Dis-
986 tinct linear and non-linear trajectories of reward and punishment reversal
987 learning during development: Relevance for dopamine's role in adolescent
988 decision making. *Developmental Cognitive Neuroscience*, 1(4), 578–590.
989 <https://doi.org/10.1016/j.dcn.2011.06.007>
- 990 Watanabe, S. (2013). A Widely Applicable Bayesian Information Criterion. *Jour-*
991 *nal of Machine Learning Research*, 14(Mar), 867–897. Retrieved October
992 30, 2019, from <http://www.jmlr.org/papers/v14/watanabe13a.html>
- 993 Wilson, R. C., & Collins, A. (2019). Ten simple rules for the computational mod-
994 eling of behavioral data. *arxiv*. <https://doi.org/10.31234/osf.io/46mbn>
- 995 Wimmer, G. E., Braun, E. K., Daw, N. D., & Shohamy, D. (2014). Episodic mem-
996 ory encoding interferes with reward learning and decreases striatal pre-
997 diction errors. *The Journal of Neuroscience: The Official Journal of the*
998 *Society for Neuroscience*, 34(45), 14901–14912. [https://doi.org/10.1523/](https://doi.org/10.1523/JNEUROSCI.0204-14.2014)
999 JNEUROSCI.0204-14.2014
- 1000 Xia, L., Master, S., Eckstein, M., Wilbrecht, L., & Collins, A. G. (2020). Learn-
1001 ing under uncertainty changes during adolescence, In *Proceedings of the*
1002 *Cognitive Science Society*.
- 1003 Yagishita, S., Hayashi-Takagi, A., Ellis-Davies, G. C. R., Urakubo, H., Ishii, S.,
1004 & Kasai, H. (2014). A critical time window for dopamine actions on the
1005 structural plasticity of dendritic spines. *Science (New York, N.Y.)*, 345(6204),
1006 1616–1620. <https://doi.org/10.1126/science.1255514>