

1 Longitudinal Changes in Value-based Learning in Middle Childhood: Distinct Contributions  
2 of Hippocampus and Striatum

3 Johannes Falck<sup>1</sup>, Lei Zhang<sup>2,3,4</sup>, Laurel Raffington<sup>5</sup>, Johannes J. Mohn<sup>6,7</sup>, Jochen Triesch<sup>8</sup>, Christine  
4 Heim<sup>6,9</sup> & Yee Lee Shing<sup>1</sup>

5  
6 <sup>1</sup>*Department of Psychology, Goethe University Frankfurt, 60629 Frankfurt am Main, Germany*

7 <sup>2</sup>*Social, Cognitive and Affective Neuroscience Unit, Department of Cognition, Emotion, and Methods  
8 in Psychology, Faculty of Psychology, University of Vienna, 1010 Vienna, Austria*

9 <sup>3</sup>*Centre for Human Brain Health, School of Psychology, University of Birmingham, Birmingham B15  
10 2TT, UK*

11 <sup>4</sup>*Institute for Mental Health, School of Psychology, University of Birmingham, Birmingham B15 2TT,  
12 UK*

13 <sup>5</sup>*Center for Lifespan Psychology, Max Planck Institute for Human Development, 14195 Berlin,  
14 Germany*

15 <sup>6</sup>*Charité – Universitätsmedizin Berlin, Institute of Medical Psychology, 10117 Berlin, Germany*

16 <sup>7</sup>*Max Planck School of Cognition, Max Planck Institute for Human Cognitive and Brain Sciences, 04103  
17 Leipzig, Germany*

18 <sup>8</sup>*Frankfurt Institute for Advanced Studies (FIAS), 60439 Frankfurt am Main, Germany*

19 <sup>9</sup>*Center for Safe & Healthy Children, The Pennsylvania State University, State College, PA 16802, USA*

20

21

22

Abstract

23 The hippocampal-dependent memory system and striatal-dependent memory system modulate  
24 reinforcement learning depending on feedback timing in adults, but their contributions during  
25 development remain unclear. In a 2-year longitudinal study, 6-to-7-year-old children performed a  
26 reinforcement learning task in which they received feedback immediately or with a short delay following  
27 their response. Children's learning was found to be sensitive to feedback timing modulations in their  
28 reaction time and inverse temperature parameter, which quantifies value-guided decision-making. They  
29 showed longitudinal improvements towards more optimal value-based learning, and their hippocampal  
30 volume showed protracted maturation. Better delayed model-derived learning covaried with larger  
31 hippocampal volume longitudinally, in line with the adult literature. In contrast, a larger striatal volume  
32 in children was associated with both better immediate and delayed model-derived learning  
33 longitudinally. These findings show, for the first time, an early hippocampal contribution to the dynamic  
34 development of reinforcement learning in middle childhood, with neurally less differentiated and more  
35 cooperative memory systems than in adults.

36

## Introduction

37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73

As children enter school during middle childhood, they must learn to act appropriately in new situations through feedback. For example, children must learn to raise their hand before speaking during class. The teacher may reinforce this behavior immediately or with a delay, which raises the question whether feedback timing modulates their learning. Here, reinforcement learning (RL)<sup>1</sup> provides a useful mechanistic framework to describe such feedback-driven value-based learning and decision-making. RL models allow to explicitly test for the influence of separate components during value-based learning, such as model-free and model-based learning<sup>2</sup>, social and non-social learning<sup>3,4</sup>, or the contribution of different memory systems<sup>5-7</sup>.

The role of feedback timing has previously been studied in relation to memory systems. The memory systems account is a theoretical framework that proposes that different types of memory are supported by distinct neural systems in the brain. Specifically, this account suggests that there are two memory systems: a hippocampal-dependent system and a striatal-dependent system. These systems modulate memory and value-based learning, and their interactive development has been of particular interest to developmental research<sup>8,9</sup>. In adults, the hippocampal-dependent memory system has been shown to contribute to episodic memory during reinforcement learning and is more engaged during feedback that is presented with a delay<sup>6,10,11</sup>, as opposed to the striatal-dependent memory system, which is more engaged after immediate feedback and supports habitual memory<sup>5,12-14</sup>. Specifically, hippocampal activation was greater during delayed feedback than during immediate feedback, whereas striatal activation was greater during immediate feedback than during delayed feedback<sup>5</sup>. The engagement of the hippocampus during delayed feedback was further supported by enhanced episodic memory for incidentally presented objects compared to objects presented with immediate feedback. Taken together, findings from adult studies suggest that feedback timing modulates the engagement of the hippocampal and striatal memory systems during value-based learning. Given the differential developmental trajectories of these systems and the impact the systems have on reinforcement learning and memory, it is important to understand whether children would show similar feedback timing modulations as previously shown in adults. In addition, whether such feedback timing modulation changes over time remains largely unexplored. To this end, in this study, we examined the contributions of hippocampal and striatal structural volumes during the longitudinal development of reinforcement learning across two years in 6-to-7-year-old children. We will introduce the key parameters in reinforcement learning and then we review the existing literature on developmental trajectories in reinforcement learning as well as on hippocampus and striatum, our two brain regions of interest.

Reinforcement learning behavior modulated by feedback timing can be modeled computationally using at least three parameters that reflect feedback-based learning and decision-making. For feedback-based learning, a learning rate parameter determines the extent to which the reward prediction error, defined as the difference between the received reward and the expected reward, influences the update of the future choice values. A higher learning rate emphasizes recent outcomes,

74 whereas a lower learning rate reflects learning integrated over a longer outcome history<sup>15</sup>. Value updates  
75 may further depend on an outcome sensitivity parameter that scales the individual magnitude of received  
76 rewards. Finally, in decision-making, the inverse temperature parameter plays a key role in determining  
77 the tendency to select the more valuable choice and quantifies choice stochasticity. A higher inverse  
78 temperature reflects more value-guided, deterministic choice behavior compared to a lower inverse  
79 temperature reflecting more random choices. Learning rates and inverse temperature have been studied  
80 extensively across development, mainly with cross-sectional studies showing mixed findings regarding  
81 their age gradients<sup>16</sup>. One study reported lower learning rates in children compared to adolescents<sup>17</sup>,  
82 while other studies found no differences<sup>18,19</sup> or even higher learning rates in children<sup>8,20</sup>. Developmental  
83 differences regarding the inverse temperature parameter are slightly more consistent, with studies  
84 reporting no differences<sup>8,21-23</sup> or higher inverse temperature with age that suggests that behavior is  
85 increasingly value-guided and less explorative<sup>17-19,24</sup>. To the best of our knowledge, outcome sensitivity  
86 has not been modeled computationally across development. However, studies that linked striatal reward  
87 activation to self-reported reward sensitivity showed increasing sensitivity from childhood to  
88 adolescence<sup>25,26</sup>.

89 In general, the inconsistencies regarding developmental differences in parameters may be due  
90 to their dependency on model and task properties<sup>27</sup>, which could be reconciled by comparing  
91 developmental changes to simulation-based optimal learning<sup>15</sup>. Such comparisons acknowledge that  
92 optimal parameter values vary depending on the context, and it has been suggested that humans develop  
93 towards more optimal parameter values from childhood into adulthood<sup>16</sup>. Importantly, to our knowledge  
94 previous reinforcement learning studies with children were cross-sectional, and only two studies  
95 investigated children under 8 years of age<sup>17,28</sup>. Cross-sectional studies, in which developmental change  
96 is inferred as a between-subject factor, do not capture the dynamics in middle childhood if individual  
97 differences are large, whereas longitudinal studies test development as a within-subject factor, which is  
98 crucial for uncovering change across time. Thus, longitudinal changes in reinforcement learning in  
99 middle childhood as well as their putative striatal and hippocampal associations remain unknown. To  
100 this end, learning rates, outcome sensitivity and inverse temperature are relevant computational  
101 parameters to study longitudinal changes in striatal and hippocampal systems during value-based  
102 learning.

103 Striatal and hippocampal contributions to reinforcement learning during middle childhood may  
104 differ as these brain regions undergo major developmental changes. Whereas earlier structural studies  
105 with relatively small sample sizes showed large developmental variability and a tendency for an earlier  
106 volume peak in the striatum than in the hippocampus<sup>29-35</sup>, a recent cross-sectional large-scale study was  
107 able to contrast striatal and hippocampal trajectories with greater granularity<sup>36</sup>. These data showed  
108 striatal volume peaks in the first decade which then declined throughout later developmental periods,  
109 whereas hippocampal volume showed a more protracted inverted-U-shaped trajectory that peaked in  
110 adolescence. Based on these structural findings, striatal and hippocampal systems are expected to

111 develop functionally at different rates<sup>37</sup>, with habit memory depending on the earlier developing striatum  
112 and episodic memory depending on the later developing hippocampus<sup>38</sup>. A direct investigation of the  
113 longitudinal development of both memory systems in childhood would shed light on whether the  
114 memory systems show a differential engagement similar to that of adults<sup>5</sup>. Such knowledge could be  
115 useful to structure learning processes according to the developmental status. For example, children's  
116 ability to learn from delayed feedback may depend on how well their hippocampus has developed. In  
117 the same study sample, we previously reported that children's hippocampal volume was related to their  
118 family's income level<sup>39</sup>. Additionally, previous research has shown that stress can reduce the  
119 effectiveness of the hippocampal-dependent memory system<sup>11</sup>. This suggests that environmental factors  
120 such as income and stress may play a role in shaping how well children learn from delayed feedback,  
121 particularly through their impact on hippocampal development. By identifying the specific  
122 environmental factors that impact children's learning and brain development, we can identify risk groups  
123 and tailor interventions to ameliorate adverse effects.

124 This study aimed to explore the development of value-based learning in children and its  
125 relationship with structural brain development over time. We hypothesized that the timing of feedback  
126 would modulate children's learning from reinforcement, and that such modulation can be captured by  
127 reinforcement learning (RL) model parameters. Additionally, we predicted that children's value-based  
128 longitudinal development would shift towards more optimal learning behavior. Regarding structural  
129 brain development, we expected the striatum to be relatively mature by middle childhood compared to  
130 the protracted hippocampal maturation. Our second objective was to investigate the relationship between  
131 value-based learning and structural brain development using longitudinal structural equation modeling.  
132 We anticipated that there would be differentiated brain-cognition links between brain volume and value-  
133 based learning. Specifically, we predicted that immediate feedback learning would be more strongly  
134 associated with striatal volume, whereas hippocampal volume would be more closely linked to delayed  
135 feedback and the facilitation of episodic memory encoding. Finally, we examined how these brain-  
136 cognition dynamics would change over time by analyzing their longitudinal changes.

137

138

## Method

139

### 140 Participants

141 Children and their parents took part in 2 waves of data collection with an interval of about 2 years (*mean*  
142 = 2.07, *SD* = 0.17, *range* = 1.69 – 2.68). The inclusion criteria for wave 1 were children attending first  
143 or second grade, no psychiatric or physical health disorders, at least one parent speaking fluent German,  
144 and born full-term ( $\geq 37$  weeks of gestation). At wave 1, 142 children (46% female, age *mean* = 7.19,  
145 *SD* = 0.46, *Range* = 6.07 - 7.98) and their parents or caregivers participated in the study. 140 children  
146 were included in the analysis (one child did not complete the probabilistic learning task, and another  
147 child was later excluded due to technical problems during the task). A subgroup of 90 children (49%

148 female, 100% right-handed), who was randomly selected, completed magnetic resonance imaging  
149 (MRI) scanning at wave 1, and 82 of them contributed to structural data after removing scans with  
150 excessive movement. At wave 2, 127 children (46% female, age  $mean = 9.25$ ,  $SD = 0.45$ ,  $Range = 8.30$   
151  $-10.2$ ) continued taking part in the study, while families of the remaining children were unable to be  
152 contacted or decided not to return to the study. 126 children at wave 2 completed the reinforcement  
153 learning task and were included in the analysis. All children at wave 2 were invited for MRI scanning,  
154 and 104 of them completed scanning (45% female, 92% right-handed). 99 children contributed to  
155 structural data, after removing scans with excessive movement. In total, 73 children contributed to the  
156 longitudinal MRI data and 126 children contributed to the longitudinal learning data. As previously  
157 reported for this study sample, we found no systematic bias due to wave 2 dropout<sup>39</sup>.

158

### 159 Procedure

160 The study consisted of a series of cognitive tasks tested during two behavioral sessions, including a  
161 reinforcement learning task, and one MRI session at wave 1<sup>39,40</sup>. Two years later, the children underwent  
162 one behavioral and one MRI session. MRI scanning was performed within three weeks of the behavioral  
163 task session. Each session lasted between 150 and 180 minutes and was scheduled either on weekdays  
164 between 2 p.m. and 6 p.m. or during weekends. Before participation, the parents provided written  
165 informed consent and children's verbal assent at both waves. All children were compensated with an  
166 honorarium of 8 euro per hour.

167

### 168 Measures

169 *Reinforcement learning task.* Children completed an adapted reinforcement learning task<sup>5</sup> in which they  
170 learned the preferred associations between four cues (cartoon characters) and two choices (round-shaped  
171 or square-shaped lolli) through probabilistic feedback (87.5 % contingent and 12.5 % non-contingent  
172 reward probability). In each trial, after an initial inter-trial interval of 0.5 s, a cue and its choice options  
173 were presented for up to 7 s until the child made a choice (Figure 1, choice phase). In the delay phase,  
174 we manipulated feedback timing. For two cues, the selected choice remained visible for 1 s (immediate  
175 feedback condition), whereas for the other two cue characters, it remained visible for 5 s before feedback  
176 was given (delayed feedback condition). A final feedback phase of 2 s indicated a reward by a green  
177 frame, and a punishment by a red frame. Inside each frame, a unique object picture was shown, which  
178 was incidentally encoded and irrelevant to the task. The child was instructed to pay attention to the  
179 feedback indicated by the frame color. In an initial practice phase of 32 trials, the child practiced the  
180 task with a fifth cartoon character not included in the actual task to avoid practice effects. The  
181 experimenter instructed the child to select the choice that was most likely to result in a reward. The  
182 Experimenter checked whether the child learned the more rewarded choice during practice and let it  
183 repeat the practice task otherwise to ensure understanding of the task. In the actual task, 128 trials were  
184 presented in four blocks and with small breaks in between. Cues were presented in a mixed, pseudo-

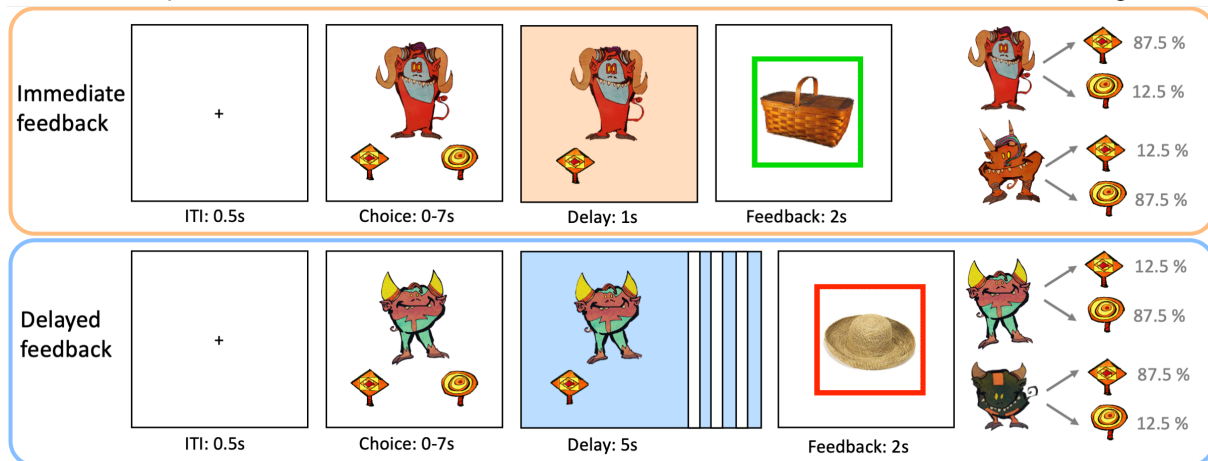
185 randomized order. A total of 64 unique objects were shown in the feedback phase, each one twice within  
186 the same feedback condition. In both delay phases, contingent choice and choice location remained the  
187 same for each cue within the task, but were balanced across participants by using four different task  
188 versions. At wave 2, four new cues replaced the previous ones to rule out memory effects.

189 *Object recognition test.* At wave 1, children were additionally tested for recognition memory on the  
190 object pictures that were incidentally encoded during reinforcement learning. A total of 80 objects (48  
191 old objects and 32 new objects) were presented in randomized order. The 48 old objects (24 for each  
192 feedback condition) were selected from the 64 old objects shown during learning based on two lists to  
193 balance the shown and omitted old objects across task versions. Each old object was shown twice during  
194 learning, but if the child failed to respond during learning, no feedback or object was shown in the trial,  
195 so some objects only appeared once. These objects were excluded at the individual level (individually  
196 missing object  $mean = 2.71$ ). At recognition, children had 4 response options ('old sure', 'old unsure',  
197 'new unsure', 'new sure') with up to 7 s to respond. The children answered verbally, and the  
198 experimenter entered their response. At wave 2, this test was excluded due to time constraints.

199

A. Two example trials

B. Reward contingencies



200

201 Figure 1. (A) Depiction of two example trials of immediate and delayed feedback conditions presented  
202 at wave 1. For immediate feedback (top panel), between choice response and feedback, cue and choice  
203 were presented for 1 s. At feedback, a green frame around the incidentally encoded object indicated a  
204 positive outcome, which appeared in 87.5% of the trials when selecting the sward-shaped lolli for this  
205 example cue. For delayed feedback (bottom panel), the delay phase between choice response and  
206 feedback lasted for 5 s. The red frame around the object indicated a negative outcome and appeared in  
207 87.5% of the trials when selecting the sward-shaped lolli for this example cue. (B) For each feedback  
208 condition, two action-outcome contingencies were learned to balance a potential choice bias. With the  
209 four task versions, the cues and outcome contingencies were counterbalanced across participants.

210

211 *Brain volume.* We extracted the bilateral brain volumes for our regions of interest, which were striatum  
212 and hippocampus. The striatum regions included nucleus accumbens, caudate and putamen. For our



213 imaging data, structural MRI images were acquired on a Siemens Magnetom TrioTim syngo 3 Tesla  
214 scanner with a 12-channel head coil (Siemens Medical AG, Erlangen, Germany) using a 3D T1-  
215 weighted Magnetization Prepared Rapid Gradient Echo (MPRAGE) sequence, with the following  
216 parameters: 192 slices; field of view = 256 mm, voxel size = 1 mm<sup>3</sup>, TR = 2500 ms; TE = 3.69 ms,  
217 flip angle = 7°, TI = 1100 ms. Volumetric segmentation was performed using the Freesurfer 6.0.0  
218 image analysis suite<sup>41</sup>. Previous studies suggested that software tools based on adult brain templates  
219 provide inaccurate segmentation for pediatric samples, which can be improved through the use of study-  
220 specific template brains<sup>42,43</sup>. Thus, we created two study-specific template brains (one for each wave)  
221 using Freesurfer's "make\_average\_subject" command. This pipeline utilized the default adult template  
222 brain registrations of the "recon-all-all" command to average surfaces, curvatures, and volumes from  
223 all subjects into a study-specific template brain. All subjects were then re-registered to this study-  
224 specific template brain to improve segmentation accuracy. Segmented images were manually inspected  
225 for accuracy and 8 cases at wave 1 and 5 cases at wave 2 were excluded for inaccurate or failed  
226 registration due to excessive motion.

227

## 228 Data analysis

229

230 *Behavioral learning performance.* As a first step, we calculated learning outcomes directly from the raw  
231 data, which were learning accuracy, win-stay and lose-shift behavior as well as reaction time. Learning  
232 accuracy was defined as the proportion to choose the more rewarding option, while win-stay and lose-  
233 shift refer to the proportion of staying with the previously chosen option after a reward and switching  
234 to the alternative choice after receiving a punishment, respectively. We used these outcomes as our  
235 dependent variables to examine the effect of the predictors feedback timing (immediate, delayed), wave  
236 (1, 2), wave 1 age, and sex (girls, boys), utilizing generalized linear mixed models (GLMM) with the R  
237 package lme4<sup>44</sup>. All reported models included random slopes for within-subject factors feedback timing  
238 and wave (see Supplementary Material 2 for the model structure). We systematically tested main effects  
239 and interactions between the predictors and their interaction had to statistically improve the predictive  
240 ability of the model to be included in the final reported model. All predictor variables were grand-mean-  
241 centered to interpret the interaction effects independent from other predictors.

242

243 *Reinforcement learning models.* As a next step, we used computational modeling to compare the  
244 learning models of basic heuristic strategies and value-based learning and to determine the model that  
245 could best capture children's trial-by-trial learning behavior. For heuristic strategies, we considered  
246 models that reflected a Win-stay-lose-shift (wsls) or a Win-stay (ws) strategy. Win-stay is a heuristic  
247 strategy in which the same action is repeated if it leads to a positive outcome in the previous trial, and  
248 Win-stay-lose-shift additionally switches to a different action if the previous outcome is negative. Note  
249 that these model-based outcomes are not identical to the win-stay and lose-shift behavior that were

250 calculated from the raw data. The use of such model-based measure offers the advantage in discerning  
251 the underlying hidden cognitive process with greater nuance, in contrast to classical approaches that  
252 directly use raw behavioral data. The models quantified the learning behavior for each individual  $I$  for  
253 each cue  $c$  and trial  $t$ . The heuristic models consisted of a weight  $w$  that reflected its degree in strategy  
254 use. In the case of reward  $r = 1$ ,  $w$  was equal to 1 for the chosen option (eg. choice A), and 0 for the  
255 unchosen option (e.g. choice B), thus maximizing win-stay, i.e., choosing A at the subsequent trial  $t + 1$ :

$$256 \quad w_{i,c,t+1,A|r=1} = 1 \text{ and } w_{i,c,t+1,B|r=1} = 0 \quad (1)$$

257 For trials  $r = 0$  (applicable only to the wsls model), model weights were the opposite, maximizing lose-  
258 shift:

$$259 \quad w_{i,c,t+1,A|r=0} = 0; w_{i,c,t+1,B|r=0} = 1 \quad (2)$$

260 The initial weights for both choices were set to  $w_{i,c,t=1} = 0.5$ . The weight  $w$  then scaled the parameter  
261  $\tau\_wsls$  or  $\tau\_ws$  to estimate the individual strategy use during decision-making. The choice probabilities  
262 were calculated using the softmax function, eg., for the chosen option  $A$ :

$$263 \quad p(A) = \frac{\exp^{w_{i,c,t,A} * \tau\_wsls_i}}{\exp^{w_{i,c,t,A} * \tau\_wsls_i} + \exp^{w_{i,c,t,B} * \tau\_wsls_i}} \quad (3)$$

264 Thus, a higher probability of strategy use was reflected by a larger value of  $\tau\_wsls$  or  $\tau\_ws$ .

265 For value-based learning, we considered a Rescorla-Wagner model and several variants based on our  
266 theoretical conceptions. The baseline value-based model  $vbm_1$  updated the value  $v$  of the selected choice  
267 ( $A$  or  $B$ ) for the next trial  $t$ . This value update was determined by calculating the difference between the  
268 received reward  $r$  and the expected value  $v$  of the selected choice, which was the reward prediction error.

269 The value update was further scaled by a learning rate  $\alpha$  ( $0 < \alpha < 1$ ):

$$270 \quad v_{i,c,t+1,A} = v_{i,c,t,A} + \alpha_i(r_{i,c,t} - v_{i,c,t,A}) \quad (4)$$

271 When the outcome sensitivity parameter  $\rho$  ( $0 < \rho < 20$ ) was included, the reward was additionally  
272 scaled at the value update:

$$273 \quad v_{i,c,t+1,A} = v_{i,c,t,A} + \alpha_i(\rho_i * r_{i,c,t} - v_{i,c,t,A}) \quad (5)$$

274 The inverse temperature parameter  $\tau$  ( $0 < \tau < 20$ ) was included in the softmax function to compute  
275 choice probabilities:

$$276 \quad p(A) = \frac{\exp^{v_{i,c,t,A} * \tau_i}}{\exp^{v_{i,c,t,A} * \tau_i} + \exp^{v_{i,c,t,B} * \tau_i}} \quad (6)$$

277 Note, however, that outcome sensitivity and inverse temperature are difficult to fit simultaneously due  
278 to non-identifiability issues<sup>45</sup>. Therefore, models including the inverse temperature fixed outcome  
279 sensitivity at 1 (inverse temperature model family), assuming no individual differences in outcome  
280 sensitivity. For the outcome sensitivity model family, outcome sensitivity was freely estimated, and the  
281 inverse temperature was fixed at 1, assuming the same degree of value-based decision behavior across  
282 individuals. Even though outcome sensitivity is usually restricted to an upper bound of 2 to not inflate  
283 outcomes at value update, this configuration led to ceiling effects in outcome sensitivity and non-  
284 converging model results. Further, this issue was not resolved when we fixed the inverse temperature at  
285 the group mean of 15.47 of the winning inverse temperature family model. It may be that in children,



286 individual differences in outcome sensitivity are more pronounced, leading to more extreme values.  
287 Therefore, we decided to extend the upper bound to 20, parallel to the inverse temperature, and all our  
288 models converged with  $R_{\text{hat}} < 1.1$ . Each model family consisted of 4 model variants  $vbm_{1-4}$  ( $1\alpha1\tau$ ,  
289  $2\alpha1\tau, 1\alpha2\tau, 2\alpha2\tau$ ) and  $vbm_{5-8}$  ( $1\alpha1\rho, 2\alpha1\rho, 1\alpha2\rho, 2\alpha2\rho$ ), in which each parameter was either  
290 separated by feedback timing or kept as a single parameter across feedback conditions. Our baseline  
291 value-based model  $vbm_1$  included a single learning rate and a single inverse temperature ( $1\alpha1\tau$ ).

292

293 *Parameter estimation.* All choice data were fitted in a hierarchical Bayesian analysis using the Stan  
294 language in R<sup>46,47</sup> adopted from the hBayesDM package<sup>48</sup>. Posterior parameter distributions were  
295 estimated using Markov chain Monte Carlo (MCMC) sampling running 4 chains each with 3,000  
296 iterations, using the first half of the chain as warmup, and group-level parameters and individual-level  
297 parameters were estimated simultaneously. The hierarchical Bayesian approach provides more stable  
298 and reliable parameter estimates as opposed to point-estimation approaches like maximum likelihood  
299 estimation<sup>49</sup>. Each model fit both wave 1 and wave 2 data at once, considering the correlation structure  
300 of the same parameter across waves, to account for within-subject dependency using the Cholesky  
301 decomposition. The Cholesky decomposition used a Lewandowski-Kurowicka-Joe prior of 2, and all  
302 other group-level parameters had a prior normal distribution, Normal (0, 0.5). Non-response trials (wave  
303 1 = 2.41%, wave 2 = 0.97% on average) were excluded in advance.

304

305 *Model simulation and model-derived learning score.* To appropriately interpret the parameter results  
306 with respect to the optimal parameter combination of the winning model, we simulated 5,000,000  
307 individual datasets using 10,000 different parameter value combinations (covering the whole range of  
308 each parameter) to identify the optimal parameter combination of the winning model that was selected  
309 by model comparison. In addition, we computed the model-derived mean choice probability of the  
310 contingent, i.e., the more rewarded option, and we referred to it as the model-derived learning score.  
311 This model-derived choice probability differs from the observed empirical choice probability (i.e., the  
312 accuracy of selecting the more rewarded option), because the model-derived learning score combines  
313 the model with the data by incorporating latent information carried out by key learning parameters. Thus,  
314 the learning score captures observed behavior based on trial-by-trial latent processes predicted by value-  
315 based models. We used this as metric to interpret the fitted posterior parameters in relation to the optimal  
316 parameter combination of our probabilistic learning task.

317

318 *Model selection and validation.* We conducted a 2-step sequential procedure for the model development  
319 and model selection. As a first step, we compared model evidence for the baseline value-based model  
320 that does not separate learning rate and inverse temperature by feedback timing ( $vbm_1:1\alpha, 1\tau$ ) to the  
321 non-value-based, heuristic strategy models that reflect Win-stay or Win-stay-lose-shift strategy behavior  
322 ( $ws, wsls$ ). As a second step, we compared model evidence for 8 value-based model variants, 4 of the

323 model family with learning rate and inverse temperature ( $1\alpha1\tau, 2\alpha1\tau, 1\alpha2\tau, 2\alpha2\tau$ ) and 4 of the model  
324 family with learning rate and outcome sensitivity ( $1\alpha1\rho, 2\alpha1\rho, 1\alpha2\rho, 2\alpha2\rho$ ). This allowed us to  
325 compare whether children showed separable effects of feedback timing on one of the model parameters.  
326 We compared the model fit using Bayesian leave-one-out cross-validation and obtained the expected  
327 log pointwise predictive density ( $elpd_{loo}$ ) using the R package `loo`<sup>50</sup>. We further computed the model  
328 weights (*Pseudo-BMA+*) using Pseudo Bayesian model averaging stabilized by Bayesian bootstrap with  
329 100,000 iterations<sup>51</sup>. To validate our models, we estimated predictive accuracy by comparing one-step-  
330 ahead model predictions with the choice data<sup>15,52</sup>. We performed parameter recovery for the winning  
331 model and model recovery by comparing it to a set of models used during model comparison  
332 (Supplementary Material 1)<sup>53</sup>.

333

#### 334 *Episodic memory at wave 1*

335 We predicted the individual corrected recognition memory (hits-false alarms) by feedback condition in  
336 a linear mixed effects model using the R package `lme4`<sup>44</sup>. Only confident (“sure”) ratings were included  
337 in the analysis, which were 98.1 % of all given responses. A total of 140 children completed the  
338 recognition memory test and 138 were included in the analysis, with two being excluded due to negative  
339 corrected recognition memory value (i.e., poor recognition memory). Age and sex were controlled for  
340 as covariates.

341

#### 342 *Longitudinal brain-cognition links*

343 We used latent change score (LCS) models to examine the longitudinal relationships between brain and  
344 learning score measures. LCS models are longitudinal structural equation models that have been widely  
345 applied to estimate developmental changes and coupling effects across domains such as the brain and  
346 cognition<sup>54,55</sup>. LCS models allow the definition of specific paths between multiple variables to test  
347 explicit hypotheses and estimate latent change from the observed variables that account for measurement  
348 error and increase testing power<sup>56</sup>. We compiled univariate LCS models for each variable separately  
349 (learning scores and brain volumes) to examine whether there was significant individual variance and  
350 change, which could be related within a multivariate LCS model as a next step. Model fit had to be at  
351 least acceptable, with a comparative fit index (*CFI*) > 0.95, standardized root mean square residual  
352 (*SRMR*) < .08 and root mean square error of approximation (*RMSEA*) < .08<sup>57</sup>. Age and sex were included  
353 as covariates at wave 1, as well as the estimated total intracranial volume (eTIV) when brain volume  
354 was included in the model. Multivariate LCS models allow to estimate meaningful brain-cognition  
355 relationships: a wave 1 covariance between brain and cognition, brain predicting change onto cognition,  
356 or vice versa, and a covariance in both brain and cognition change scores (wave 1 to wave 2). Before  
357 compiling the variables into an LCS model, they were checked for outliers  $\pm 4$  *SD* around the mean. We  
358 identified one outlier for the learning rate at wave 2, which was removed for the explorative LCS model  
359 that included model parameters. There were no further outliers in other cognitive variables or brain

360 volumes. Continuous variables were standardized to the wave 1 measure so that wave 2 values represent  
 361 the change from wave 1, sex was contrast-coded (girls = 1, boys = -1).

362

363

364

## Results

365

366 Behavioral results

367

368 First, we were interested in whether children showed behavioral differences between waves and  
 369 feedback timing. A descriptive overview is provided in Table 1 and Figure 2. The details of the reported  
 370 GLMM models, including the random effects structure and the effects of age and sex, are described in  
 371 the Supplementary Material 2. Since some children were poor learners who failed to reach 50 % average  
 372 accuracy in their last 20 trials (13 children at wave 1 and 6 children at wave 2), we also performed  
 373 behavioral analyses with a reduced dataset in which results remained unchanged (Supplementary  
 374 Materials 6).

375

376 Table 1. Descriptive behavioral results of dependent variables Accuracy (ACC, probability correct),  
 377 win-stay probability (WS), lose-shift probability (LS), and reaction time (RT, in seconds), as well as  
 378 mixed model fixed effects that predicted these dependent variables.

	Descriptive Results				Mixed Model Effects	
	Wave 1		Wave 2		Wave	Feedback
	Ime	Del	Ime	Del		
ACC	0.69 (0.46)	0.70 (0.46)	0.79 (0.41)	0.80 (0.40)	↑ W2	–
WS	0.81 (0.39)	0.80 (0.40)	0.88 (0.32)	0.88 (0.32)	↑ W2	–
LS	0.47 (0.50)	0.50 (0.50)	0.42 (0.49)	0.42 (0.49)	↓ W2	–
RT	2.10 (1.31)	2.07 (1.29)	1.70 (1.02)	1.67 (1.00)	↓ W2	↓ Del

379 *Note.* Mean (standard deviation) of the variables, split by wave and feedback timing, is reported in the  
 380 table. Mixed model effects and their directionality (increasing ↑ or decreasing ↓) predicting the  
 381 dependent variables. W2 = Wave 2, Ime = Immediate feedback, Del = Delayed feedback.

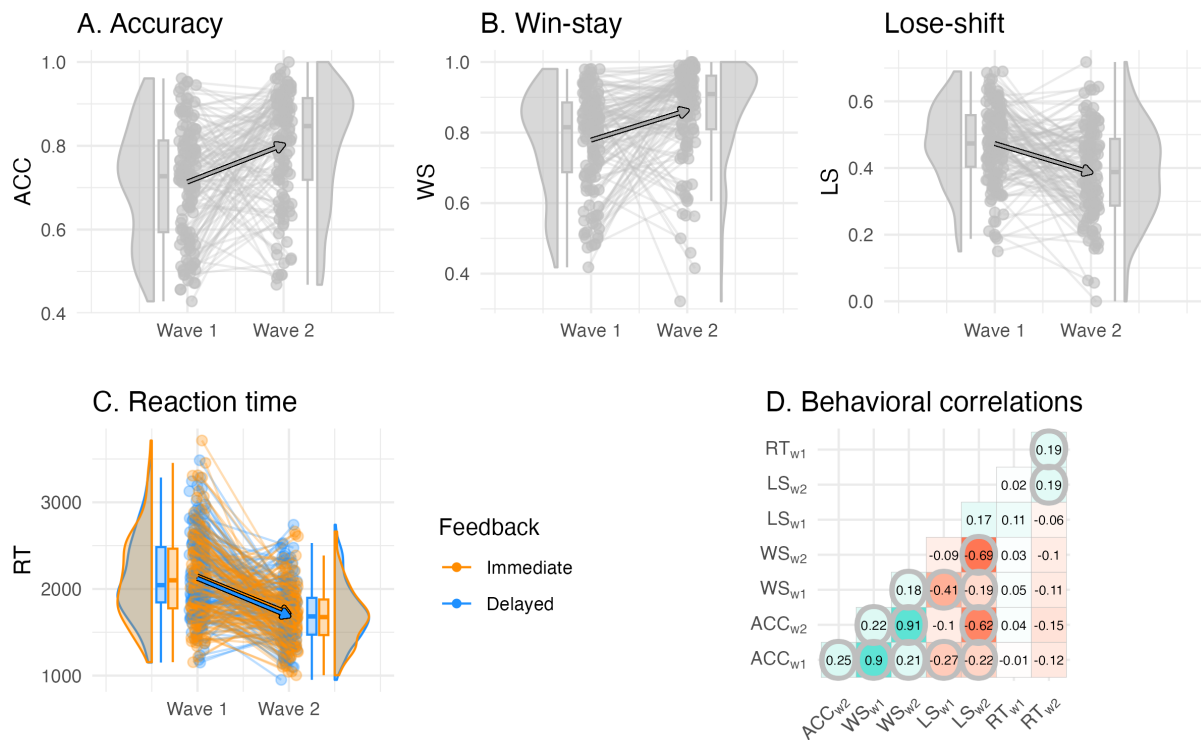
382

383 *Children's learning improved between waves.* With the complete dataset, we found that increased  
 384 learning accuracy (i.e., the probability of choosing the more rewarding option) was predicted at wave 2  
 385 compared to wave 1, but there were no differences in accuracy by feedback timing ( $\beta_{wave=2} = .550$ ,  $SE$   
 386  $= .061$ ,  $z = 8.97$ ,  $p < .001$ ,  $\beta_{feedback=delayed} = .013$ ,  $SE = .024$ ,  $z = 0.54$ ,  $p = .590$ ). Furthermore, win-  
 387 stay probability increased and lose-shift probability decreased longitudinally, again without differences  
 388 by feedback timing (WS:  $\beta_{wave=2} = .586$ ,  $SE = .071$ ,  $z = 8.22$ ,  $p < .001$ ,  $\beta_{feedback=delayed} = .023$ ,  $SE$

## Longitudinal Changes in Value-based Learning in Middle Childhood

12

389 = .033,  $z = 0.69$ ,  $p = .489$ ; LS:  $\beta_{wave=2} = -.252$ ,  $SE = .037$ ,  $z = -6.87$ ,  $p < .001$ ,  $\beta_{feedback=delayed} = .030$ ,  
 390  $SE = .022$ ,  $z = 1.37$ ,  $p = .169$ ). Reaction times were faster at wave 2 compared to wave 1, and they were  
 391 faster for delayed compared to immediate feedback trials ( $\beta_{wave=2} = -.221$ ,  $SE = 22.8$ ,  $t(df_{Satterthwaite} =$   
 392  $135) = -9.70$ ,  $p < .001$ ,  $\beta_{feedback=delayed} = -13.8$ ,  $SE = 6.59$ ,  $t(df_{Satterthwaite} = 136) = -2.10$ ,  $p = .038$ ). To  
 393 summarize, children's average accuracy improved over 2 years, while their win-stay probability  
 394 increased and their lose-shift probability decreased between waves. Children were able to respond faster  
 395 to cues paired with delayed feedback compared to cues paired with immediate feedback, and they  
 396 became faster in their decision-making across waves (see mixed model effects overview in Table 1). Of  
 397 note, reaction times were largely uncorrelated with accuracy and switching behavior (win-stay, lose-  
 398 shift), while accuracy and switching behavior showed significant correlations at both waves (Figure 2D).  
 399  
 400



401  
 402 Figure 2. Individual differences in the behavioral reinforcement learning outcomes and their longitudinal  
 403 change. (A) Accuracy did not differ by feedback timing and increased between waves. (B) Win-stay and  
 404 lose-shift proportion did not differ by feedback timing, and win-stay increased and lose-shift proportion  
 405 decreased between waves. (C) Reaction time differed by feedback timing, in which decisions for cues  
 406 learned with delayed feedback were faster, and reaction times were faster at wave 2 compared to wave  
 407 1. (D) Correlations between behavioral outcomes reveal that learning accuracy was primarily correlated  
 408 with the win-stay and lose-shift probabilities both within and between waves, but was uncorrelated to  
 409 reaction time. Significant correlations are circled,  $p$ -values were adjusted for multiple comparisons using  
 410 bonferroni correction.  
 411

412 Modeling results

413

414 *Children's behavior was best described by value-based learning.* We conducted a 2-step sequential  
 415 procedure for model development and model selection. Model comparison using leave-one-out cross  
 416 validation showed evidence in favor of the value-based learning model, reflected in the highest expected  
 417 log pointwise predictive density and highest model weights, confirming that children's learning  
 418 behavior in the longitudinal data can generally be better described by a value-based rather than by a  
 419 heuristic strategy model ( $\Delta elpd_{loo} = -15154.9$ ,  $pseudo-BMA+ = 1$ , Table 2). Children whose individual  
 420 fit was better for a heuristic model ( $wsls$ ) than for the value-based model ( $vbm_1$ ), were at both waves  
 421 more likely to be poor learners (defined as an accuracy below 50% in the last 20 trials). Taken together,  
 422 children's learning behavior was best described by a value-based model, and a heuristic strategy model  
 423 captured more poor learners compared to a value-based model.

424

425 Table 2. Model comparison results.

Model	Parameters	$\Delta elpd_{loo}$ [SE]	$\Sigma elpd_{loo}$ [mean]	$pseudo-BMA+$
Step 1: heuristic strategy models and value-based learning model				
$vbm_1$	$1\alpha, 1\tau$	0 [0]	-15154.9 [-0.45]	1
$ws$	$1\tau_{ws}$	-1327.7 [159.5]	-16482.7 [-0.49]	<0.01
$wsls$	$1\tau_{wsls}$	-4247.3 [284.8]	-19402.3 [-0.58]	0
Step 2: value-based learning models				
<b><math>vbm_3</math></b>	<b><math>1\alpha, 2\tau</math></b>	<b>0 [0]</b>	<b>-15045.3 [-0.45]</b>	<b>0.73</b>
$vbm_7$	$1\alpha, 2\rho$	-2.93 [2.92]	-15048.2 [-0.45]	0.24
$vbm_6$	$2\alpha, 1\rho$	-24.34 [8.85]	-15069.6 [-0.45]	<0.01
$vbm_8$	$2\alpha, 2\rho$	-29.71 [15.95]	-15075.0 [-0.45]	0.02
$vbm_4$	$2\alpha, 2\tau$	-43.34 [14.89]	-15088.6 [-0.45]	<0.01
$vbm_2$	$2\alpha, 1\tau$	-46.45 [13.97]	-15091.7 [-0.45]	<0.01
$vbm_5$	$1\alpha, 1\rho$	-59.01 [7.59]	-15104.3 [-0.45]	<0.01
$vbm_1$	$1\alpha, 1\tau$	-109.63 [11.98]	-15154.9 [-0.45]	<0.01

426 *Note.* Model = heuristic ( $ws$ ,  $wsls$ ) and value-based models ( $vbm_{1-8}$ ) that were compared against each  
 427 other. Parameters = corresponding model parameters learning rate  $\alpha$ , inverse temperature  $\tau$  and  
 428 outcome sensitivity  $\rho$ .  $\Delta elpd_{loo}[SE]$  = difference in the Bayesian leave-one-out cross-validation  
 429 estimate of the expected log pointwise predictive density relative to the winning model and its standard  
 430 errors.  $\Sigma elpd_{loo}[mean]$  = sum of expected log pointwise predictive density of all 33,460 trials,  
 431 including all participants and waves, and trial mean.  $Pseudo-BMA+$  = model weight for relative model  
 432 evidence using Bayesian model averaging stabilized by Bayesian bootstrap using 100,000 iterations.

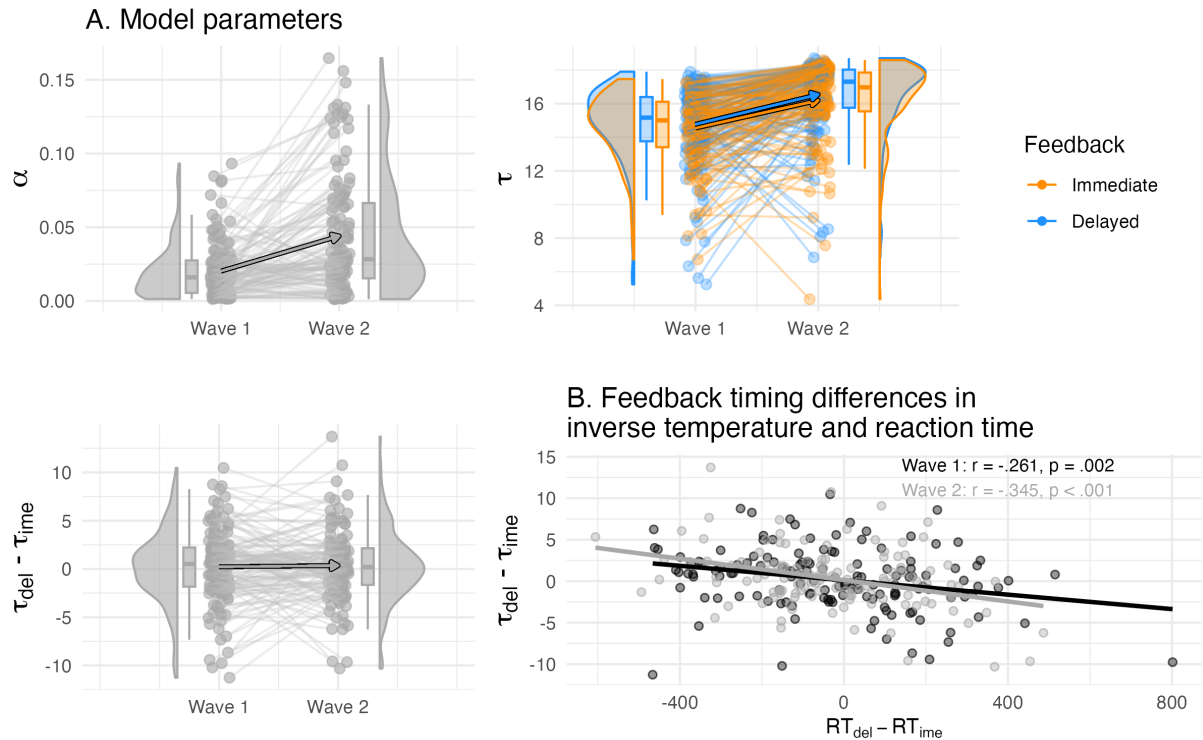
433

434 *Feedback timing modulated choice stochasticity.* Model  $vbm_3 (1\alpha 2\tau)$  showed the largest model  
435 evidence, reflected in the highest expected log pointwise predictive density and highest model weights  
436 and suggests that feedback timing affected the inverse temperature, but not the learning rate or outcome  
437 sensitivity ( $elpd_{loo} = -15045.3$ ,  $pseudo-BMA^+ = 0.73$ , Table 2). Table 3 and Figure 3A provide a  
438 descriptive overview of the winning model parameters. Of note, there were only small differences in  
439 model fit ( $elpd_{loo}$ ) to the second-best model ( $vbm_7, 1\alpha 2\rho$ ,  $\Delta elpd_{loo} = -2.93$ ,  $elpd\_SE_{loo} = 2.92$ ,  
440  $pseudo-BMA^+ = 0.24$ ), which suggests a potential separable feedback timing effect on outcome  
441 sensitivity. We also performed the model comparison with a reduced dataset in which the winning model  
442 remained the same (Supplementary Materials 6). The average inverse temperature did not differ by  
443 feedback condition, but showed large within-person condition differences at both waves, indicating  
444 individual differences in feedback timing modulation (wave 1:  $\Delta\tau_{del-ime}$   $Mean = 0.22$ ,  $SD = 3.80$ ,  
445  $Range = 21.74$ , wave 2:  $\Delta\tau_{del-ime}$   $Mean = 0.35$ ,  $SD = 3.70$ ,  $Range = 24.03$ ). The correlations between  
446 the parameters are shown in Supplementary Material 3.

447 Since reaction times were predicted by feedback timing behaviorally, and inverse temperature is  
448 assumed to reflect decision-making, we were interested in whether differences in reaction time were  
449 related to inverse temperature differences. Indeed, at both waves, children who responded faster during  
450 delayed compared to immediate feedback had a higher inverse temperature at delayed compared to  
451 immediate feedback (wave 1:  $r = -.261$ ,  $t(df = 138) = -3.18$ ,  $p = .002$ , wave 2:  $r = -.345$ ,  $t(df = 124) = -$   
452  $4.10$ ,  $p < .001$ , Figure 3B). Taken together, children's learning behavior was best described by a value-  
453 based model, where feedback timing modulated individual differences in the choice rule during value-  
454 based learning. Interestingly, the differences in the choice rule and reaction time were correlated.  
455 Specifically, more value-guided choice behavior (i.e., higher inverse temperature) was related to faster  
456 responses during delayed feedback relative to immediate feedback, suggesting a link between model  
457 parameter and behavior in relation to feedback timing.

458





459

460 Figure 3. (A) Individual differences in the learning rate and inverse temperature of the winning model  
 461 and their longitudinal change. The inverse temperature  $\tau$  but not learning rate  $\alpha$  was separated by  
 462 feedback timing, and both increased between waves in their values (top panel). The condition difference  
 463 in the inverse temperature did not differ on average, but showed individual differences (bottom left  
 464 panel). (B) The condition differences in the inverse temperature correlated with reaction time, i.e., higher  
 465 delayed compared to immediate inverse temperature was related to faster delayed compared to  
 466 immediate reaction time.

467

468 Table 3. Description of model parameters from the winning value-based model  $vbm_3$ .

	Wave 1					Wave 2				
	$\alpha$	$\tau_{Ime}$	$\tau_{Del}$	$ls_{Ime}$	$ls_{Del}$	$\alpha$	$\tau_{Ime}$	$\tau_{Del}$	$ls_{Ime}$	$ls_{Del}$
Mean	0.02	14.6	14.8	0.73	0.73	0.05	16.2	16.5	0.82	0.82
SD	0.02	2.04	2.37	0.12	0.13	0.04	2.37	2.21	0.13	0.13
Min	<0.01	6.73	5.25	0.53	0.53	<0.01	4.37	6.85	0.53	0.53
Max	0.09	17.5	17.9	0.94	0.94	0.22	18.6	18.7	0.96	0.96

469 Note.  $\alpha$  = learning rate across feedback timing,  $\tau_{Ime}/ls_{Ime}$  = inverse temperature and learning score for  
 470 immediate feedback,  $\tau_{Del}/ls_{Del}$  = inverse temperature and learning score for delayed feedback.

471

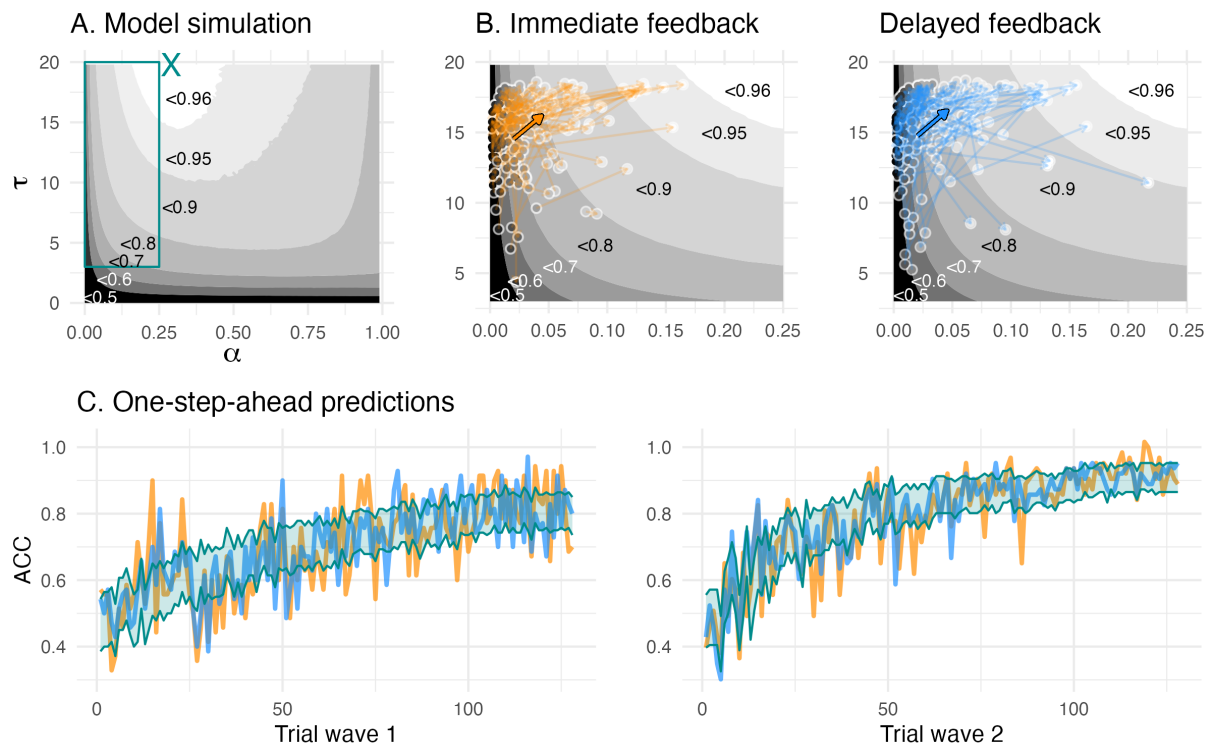
472 Children's value-based learning became more optimal. Next, we compared the parameter space  
 473 according to model simulation (Figure 4A) with the empirical posterior parameters fitted by the winning  
 474 model (Table 3, Figure 4B) to determine whether children increased their value-based learning towards  
 475 more optimal parameter combinations. Both fitted and simulated parameter combinations allowed us to  
 476 derive a learning score that captured learning performance according to the winning value-based model.

477 Note that the learning score was defined as the average choice probability for the more rewarded choice  
478 option. We refer to these model-derived choice probabilities as learning score, since they reflect value-  
479 based learning and combine information of learned values, that depend on the learning rate, and values  
480 translated into choice probabilities, that depend on the inverse temperature. Thus, a higher learning score  
481 reflects more optimal value-based learning. We simulated 10,000 parameter combinations and created  
482 a learning score map according to each parameter combination (Figure 4A). The optimal parameter  
483 combination was at a learning rate  $\alpha = 0.29$ , and an inverse temperature  $\tau = 19.8$ , and with an average  
484 learning score of 96.5 % (Figure 4A). Children's fitted learning rates ranged 0.01 – 0.22 and inverse  
485 temperature 6.73 – 18.70 and were outside the parameter space of a learning score above 96 % (Table  
486 3 and Figure 4A). The average longitudinal increases in learning rate and inverse temperature were  
487 mirrored by average increases in the learning scores, confirming our prediction that their parameters  
488 developed towards optimal value-based learning (arrow in Figure 4B). We further found that the average  
489 longitudinal change in win-stay and lose-shift proportion also developed towards more optimal value-  
490 based learning (Supplementary Material 4).

491

492 *Model validation.* To validate our winning model  $vbm_3$ , we estimated its predictive accuracy by  
493 comparing one-step-ahead model predictions with the choice data. The one-step ahead predictions of  
494 the winning model captured children's choices overall well, with predictive accuracies of 65.3 % at  
495 wave 1 and 75.7 % at wave 2 (Figure 4C). Further, our winning model showed a good parameter  
496 recovery for learning rate ( $r = 0.85$ ) and inverse temperature ( $r = 0.75 - 0.77$ ). Our winning model  
497 showed excellent on the group level (100%) when comparing it to a set of models used during model  
498 comparison ( $vbm_1, vbm_7, wsls$ ). The individual model recovery was lower (58%), with 35% of the  
499 simulated winning model fitting best on our baseline model  $vbm_1$  with a single inverse temperature,  
500 which likely reflects the noisy property of the inverse temperature (Supplementary Material 1).

501



502

503 Figure 4. (A) The model simulation depicts parameter combinations and simulation-based average  
504 learning scores. The cyan “X” in the middle top depicts the optimal parameter combination where  
505 average learning scores were at 96.5 %, and the cyan rectangle depicts the space of the fitted parameter  
506 combinations, (B) Enlarged view of the space of fitted parameter combinations. The colored arrows  
507 depict mean change (bold arrow) and individual change (transparent arrows) of the fitted parameters.  
508 The greyscale gradient-filled dots, that are connected by the arrows, depict the individual learning score,  
509 while the the greyscale gradient in the background depicts the simulated average learning score. The  
510 mean change reveals an overall change towards the higher, i.e., more optimal, learning scores. (C) One-  
511 step-ahead posterior predictions of the winning model for each wave. The colored lines depict averaged  
512 trial-by-trial task behavior for each feedback condition, and a cyan ribbon indicates the 95% highest  
513 density interval of the one-step-ahead prediction using the entire posterior distribution.

514

515 Longitudinal brain-cognition links

516

517 *Significant longitudinal change in brain and cognition.* We first performed univariate LCS model  
518 analyses to estimate a latent change score of immediate and delayed learning scores as well as striatal  
519 and hippocampal volumes (see descriptive changes in Figure 5B-C). All four variables of interest  
520 showed significant positive mean changes and variances, and all univariate models provided a good fit  
521 to the data (Supplementary Material 5). This allowed us to further relate the differences in structural  
522 brain changes to changes in learning.

523

524 *Hippocampal volume exhibited more protracted development during middle childhood.* We next fitted  
525 a bivariate LCS model to compare striatal and hippocampal change scores. We theorized that by middle  
526 childhood, the striatum would be relatively mature, whereas the hippocampus continues to develop. We  
527 progressively constructed multiple LCS models to test this idea. First, the bivariate LCS model provided  
528 a good data fit ( $\chi^2(14) = 10.09$ ,  $CFI = 1.00$ ,  $RMSEA(CI) = 0(0-.06)$ ,  $SRMR = .04$ ). We then further  
529 fitted two constrained models, to see whether setting the mean striatal change or the mean hippocampal  
530 change to 0 would lead to a drop in the model fit. Compared to the unrestricted model, the constrained  
531 model that assumed no striatal change did not lead to a drop in model fit ( $\Delta\chi^2(1) = 2.74$ ,  $p = .098$ ),  
532 whereas the model that assumed hippocampal change dropped in model fit ( $\Delta\chi^2(1) = 12.69$ ,  $p < .001$ ).  
533 Finally, we tested a more stringent assumption of equal change for striatal and hippocampal volumes,  
534 in which the model dropped in model fit compared to the unrestricted model ( $\Delta\chi^2(1) = 18.04$ ,  $p < .001$ )  
535 and suggests that striatal and hippocampal change differed. Together, these results support our  
536 postulation of separable maturational brain trajectories in our study sample, suggesting that the  
537 hippocampus continued to grow in middle childhood, whereas striatal volume increased less.

538

539 *Hippocampal and striatal volume showed distinct associations to learning.* We fitted a four-variate LCS  
540 model to test our prediction of selective brain-cognition links. Specifically, we assumed a larger  
541 contribution of striatal volume at immediate learning, and a larger contribution of hippocampal volume  
542 at delayed learning. The LCS model provided good data fit ( $\chi^2(27) = 15.4$ ,  $CFI = 1.00$ ,  $RMSEA(CI) =$   
543  $0(0 - .010)$ ,  $SRMR = .045$ ), and all relevant paths are shown in Figure 5D (see Table 4 for a detailed  
544 model overview). For the striatal associations to cognition, we found that wave 1 striatal volume  
545 covaried with both immediate learning score and delayed learning score ( $\phi_{STR_{w1}, LS_{i,w1}} = 0.19$ ,  $z = 2.52$ ,  
546  $SE = 0.07$ ,  $p = .012$ ,  $\phi_{STR_{w1}, LS_{d,w1}} = 0.18$ ,  $z = 2.37$ ,  $SE = 0.07$ ,  $p = .018$ ). Constraining the striatal  
547 association to immediate learning to 0 worsened the model fit relative to the unrestricted model ( $\Delta\chi^2(1)$   
548  $= 5.66$ ,  $p = .017$ ), which was the same when constraining the striatal association to delayed learning to  
549 0 ( $\Delta\chi^2(1) = 5.14$ ,  $p = .023$ ). In summary, larger striatal volume was associated with better learning  
550 scores for both immediate and better delayed feedback. This pattern remained the same in the results of  
551 the reduced dataset (Supplementary Material 6).

552 Hippocampal volume, on the other hand, only covaried with delayed learning at wave 1 ( $\phi_{HPC_{w1}, LS_{d,w1}} =$   
553  $0.14$ ,  $z = 2.05$ ,  $SE = 0.07$ ,  $p = .041$ ), not with immediate learning score ( $\phi_{HPC_{w1}, LS_{i,w1}} = 0.12$ ,  $z = 1.68$ ,  
554  $SE = 0.07$ ,  $p = .092$ ). Fixing the path between hippocampal volume and delayed learning to 0 worsened  
555 the model fit relative to the unrestricted model ( $\Delta\chi^2(1) = 4.19$ ,  $p = .041$ ), but not when its path to  
556 immediate learning was constrained to 0 ( $\Delta\chi^2(1) = 2.94$ ,  $p = .086$ ). This suggests that larger hippocampal  
557 volume was specifically associated with better delayed learning. In the results of the reduced dataset,  
558 the hippocampal association to the delayed learning score was no longer significant, suggesting a  
559 weakened pattern when excluding poor learners (Supplementary Material 6). It is likely that the  
560 exclusion reduced the group variance for hippocampal volume and delayed learning score in the model.

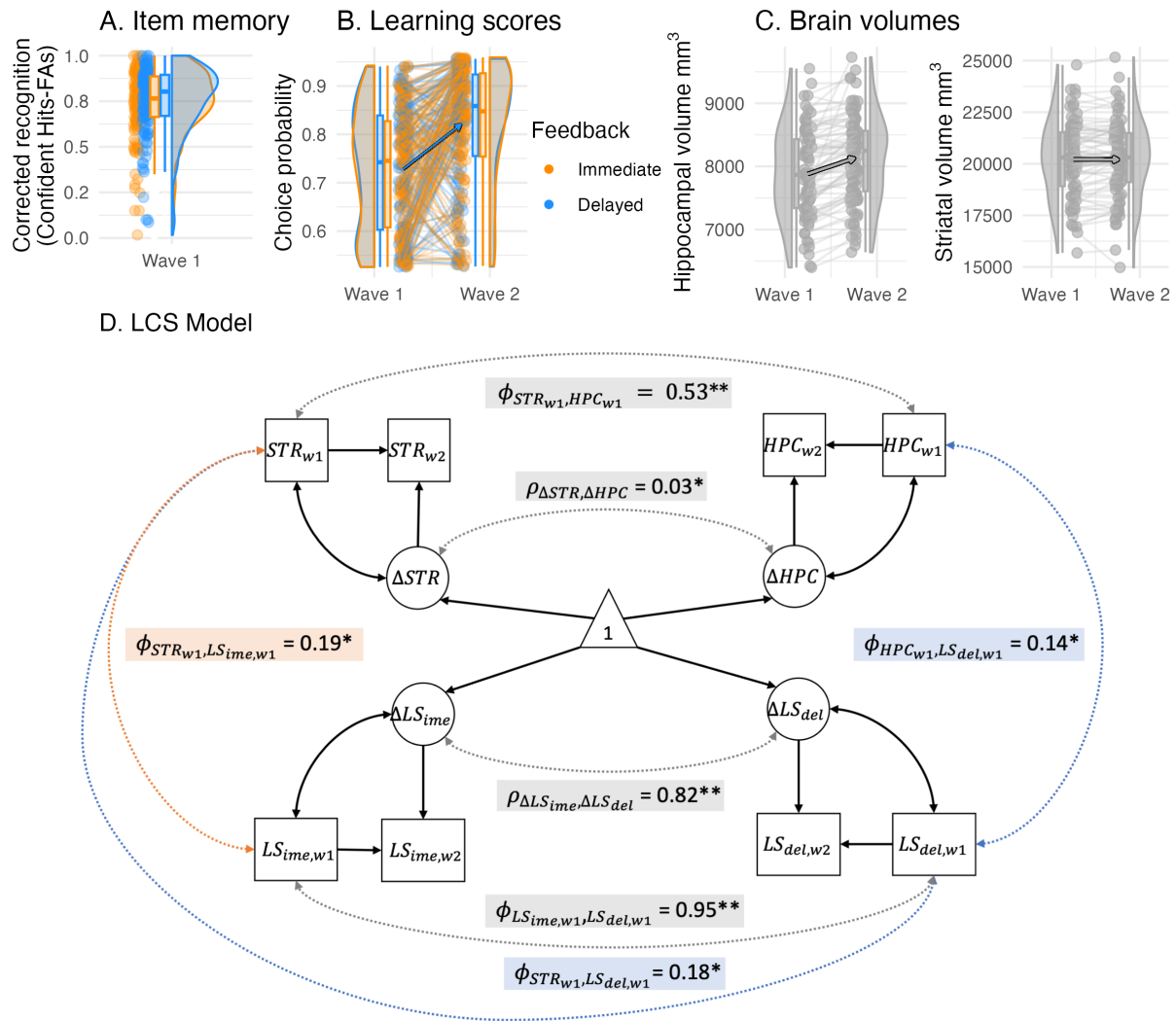
561 As a next step, the associations between striatum and hippocampus to immediate or delayed learning  
562 was directly compared against each other. A model equal-constraining striatal and hippocampal paths  
563 to immediate learning ( $\Delta\chi^2(1) = 0.41, p = .521$ ) and another model equal-constraining these paths to  
564 delayed learning ( $\Delta\chi^2(1) = 0.14, p = .707$ ) did not lead to a worse model fit compared to the unrestricted  
565 model, which suggests that the brain-cognition links have considerable overlap. This is in line with the  
566 high wave 1 covariance and change-change covariance within the brain and cognition domain (see Table  
567 4). We found no longitudinal links between the brain and cognition domains, which suggests that the  
568 found brain-cognition links at wave 1 remained longitudinally stable (see Supplementary Material 5 for  
569 an exploratory LCS model that related the model parameters to striatal and hippocampal volume).  
570 Taken together, the confirmatory LCS model results were in line with our predictions of a relatively  
571 larger involvement of the hippocampus during delayed feedback learning, but the findings on striatal  
572 volume disconfirmed a selective association with immediate feedback learning and suggest a more  
573 general role of the striatum in both learning conditions.

574

575 *No evidence for enhanced episodic memory during delayed feedback.* Finally, we investigated whether  
576 a hippocampal contribution at delayed feedback would selectively enhance episodic memory. Episodic  
577 memory, as measured by individual corrected object recognition memory (hits - false alarms) of  
578 confident (“sure”) ratings, showed at trend better memory for items shown in the delayed feedback  
579 condition ( $\beta_{feedback=delayed} = .009, SE = .005, t(df = 137) = 1.80, p = .074$ , see Figure 5A). Note that  
580 in the reduced dataset, delayed feedback predicted enhanced item memory significantly (Supplementary  
581 Material 6). The inclusion of poor learners in the complete dataset may have weakened this effect because  
582 their hippocampal function was worse and was not involved in learning (nor encoding), regardless of  
583 feedback timing. To summarize, there was inconclusive support for enhanced episodic memory during  
584 delayed compared to immediate feedback, calling for future study to test the postulation of a selective  
585 association between hippocampal volume and delayed feedback learning.

586

Longitudinal Changes in Value-based Learning in Middle Childhood



587  
588 Figure 5. (A) Recognition memory (corrected recognition = hits - false alarms) for objects presented  
589 during delayed feedback was only enhanced at trend. (B) Learning scores depicted here were used in  
590 the LCS analyses. Learning scores were the model-derived choice probability of the contingent choice  
591 using fitted posterior parameters. (C) Hippocampal and striatal volumes increased between waves, while  
592 hippocampal volume increased most. (D) A four-variate latent change score (LCS) model that included  
593 striatal and hippocampal volumes as well as immediate and delayed learning scores. Depicted are  
594 significant paths cross-domain (brain-cognition, dashed lines) and within-domain (brain or cognition,  
595 solid lines), other paths are omitted for visual clarity and are summarized in Table 4. Depicted brain-  
596 cognition links included  $\phi_{STR_{w1}, LS_{ime, w1}}$  (covariance between striatal volume and immediate learning  
597 score at wave 1), as well as  $\phi_{HPC_{w1}, LS_{del, w1}}$  and  $\phi_{STR_{w1}, LS_{del, w1}}$  (covariances between hippocampal and  
598 striatal volumes and delayed learning score at wave 1). Brain links included  $\phi_{STR_{w1}, HPC_{w1}}$  and  
599  $\rho_{\Delta STR, \Delta HPC}$  (wave 1 covariance and change-change covariance), and similarly, cognition links included  
600  $\phi_{LS_{ime, w1}, LS_{del, w1}}$  and  $\rho_{\Delta LS_{ime}, \Delta LS_{del}}$ . Covariates included age, sex and estimated total intracranial  
601 volume. \*\* denotes significance at  $\alpha < .001$ , \* at  $\alpha < .05$ .



603 Table 4. Parameter estimates of a four-variate latent change score model that includes brain (striatal and  
604 hippocampal volume) and cognition domains (immediate and delayed learning score)

	<i>STR</i>	<i>LS<sub>ime</sub></i>	<i>HPC</i>	<i>LS<sub>del</sub></i>
Model fit: $\chi^2 = 15.4$ , $df = 27$ , $CFI = 1$ , $RMSEA (CI) = 0 (0-0.01)$ , $SRMR = 0.045$				
Mean change $\Delta$	0.06* (0.03)	0.76** (0.08)	0.38** (0.04)	0.75** (0.08)
wave 1 variance $\sigma$	fixed to 1	fixed to 1	fixed to 1	fixed to 1
change variance $\sigma_{\Delta}$	0.07** (0.01)	0.88** (0.10)	0.18* (0.07)	0.83** (0.10)
Intercept-change regression $\beta$	-0.04 (0.04)	-0.83* (0.29)	-0.16* (0.06)	-0.73* (0.27)
Wave 1 covariates				
age onto Intercept $\phi$	0.19 (0.10)	-0.05 (0.08)	0.29* (0.08)	0.08 (0.08)
sex onto Intercept $\phi$	-0.42** (0.07)	-0.14 (0.07)	-0.47** (0.07)	-0.11 (0.07)
eTIV onto Intercept $\phi$	0.68** (0.05)	–	0.70** (0.05)	–
Brain-cognition links (cross-domain)				
STR- <i>LS<sub>ime</sub></i>				
wave 1 covariation $\phi$	<b>0.19* (0.07)</b>	<b>0.18* (0.07)</b>	0.12 (0.07)	<b>0.14* (0.07)</b>
change-change covariance $\rho$	<0.01 (0.03)	<0.01 (0.03)	-0.06 (0.05)	-0.07 (0.05)
wave 1 brain onto cognition change $\gamma$	0.25 (0.13)	0.22 (0.12)	0.05 (0.11)	0.06 (0.10)
wave 1 cognition onto brain change $\gamma$	-0.19 (0.13)	0.21 (0.13)	0.05 (0.10)	<0.01 (0.10)
Brain links (within-domain)				
STR- <i>HPC</i>				
wave 1 covariation $\phi$	<b>0.53** (0.07)</b>			
change-change covariance $\rho$	<b>0.03* (0.01)</b>			
wave 1 striatum onto hippocampal change $\gamma$	0.06 (0.05)			
wave 1 hippocampus onto striatal change $\gamma$	0.02 (0.03)			
Cognition links (within-domain)				
<i>LS<sub>ime</sub></i> - <i>LS<sub>del</sub></i>				
wave 1 covariation $\phi$	<b>0.95** (0.10)</b>			
change-change covariance $\rho$	<b>0.82** (0.10)</b>			
wave 1 <i>LS<sub>ime</sub></i> into <i>LS<sub>del</sub></i> change $\gamma$	-0.07 (0.27)			
wave 1 <i>LS<sub>del</sub></i> into <i>LS<sub>ime</sub></i> change $\gamma$	0.06 (0.28)			

605 Parameter estimates in bold are the paths of interest depicted in Figure 5D. Standard errors are shown in  
606 parentheses. eTIV = estimated total intracranial volume. \*\* denotes significance at  $\alpha < .001$ , \* at  $\alpha < .05$ . sex  
607 coded as 1 = girls, -1 = boys.

608

609

## Discussion

610

611 In this study, we examined the longitudinal development of value-based learning in middle childhood  
612 and its associations with striatal and hippocampal volumes that were predicted to differ by feedback  
613 timing. Children improved their learning in the 2-year study period. Behaviorally, learning was  
614 improved by an increase in accuracy and a reduction in reaction time (i.e., faster responses). Further,  
615 children's switching behavior improved by an increase in win-stay and a decrease in lose-shift behavior.

616 Computationally, learning was enhanced by an increase in learning rate and inverse temperature, which  
617 together constituted more optimal value-based learning. Further, feedback timing modulated specifically  
618 the inverse temperature. In terms of brain structures, we found that longitudinal changes in hippocampal  
619 volume were larger compared to striatal volume, which suggests more protracted hippocampal  
620 maturation. The brain-cognition links were longitudinally stable and partially confirmed our hypotheses.  
621 In line with previous adult literature and our assumption, hippocampal volume was more strongly  
622 associated with delayed feedback learning. Contrary to our expectations, episodic memory performance  
623 was not enhanced under delayed feedback compared to immediate feedback. Furthermore, striatal  
624 volume unexpectedly was associated with both immediate and delayed feedback learning, suggesting a  
625 common involvement of the striatum during value-based learning in middle childhood across timescales.

626

627 Children's learning improvement between waves was described behaviorally by increased win-stay and  
628 decreased lose-shift behavior. Our finding is in line with cross-sectional studies in the developmental  
629 literature that reported increased learning accuracy and win-stay behavior<sup>58,59</sup>. Our longitudinal dataset  
630 with younger children further suggests that learning change is not only accompanied by increased win-  
631 stay, but also decreased lose-shift behavior. We found lower learning performance and less optimal  
632 switching behavior in girls compared to boys, which could point to sex differences for reinforcement  
633 learning during middle childhood (Supplementary Material 2). Previous studies have found both male  
634 and female advantages depending on their age and the type of learning task<sup>38,60,61</sup>. Alternatively, sex  
635 differences may have been driven by confounding variables not included in the analysis.

636

637 Computationally, we found longitudinally increased and more optimal learning rate and inverse  
638 temperature, as shown by simulation data, that add to the growing literature of developmental  
639 reinforcement learning<sup>16</sup>. Adult studies that examined feedback timing during reinforcement learning  
640 reported average learning rates range from 0.12 to 0.34<sup>5,13,14</sup>, which are much closer to the simulated  
641 optimal learning rates of 0.29 than children's average learning rates of 0.02 and 0.05 at wave 1 and 2 in  
642 our study. Therefore, it is likely that individuals approach adult-like optimal learning rates later during  
643 adolescence. However, the differences in learning rate across studies have to be interpreted with caution.  
644 The differences in the task and the analysis approach may limit their comparability<sup>15,27</sup>. Task properties  
645 such as the trial number per condition differed across studies. Our study included 32 trials per cue in  
646 each condition, while in adult studies, the trials per condition ranged from 28 to 100<sup>5,13,14</sup>. Optimal  
647 learning rates in a stable learning environment were at around 0.25 for 10 to 30 trials<sup>15</sup>, another study  
648 reported a lower optimal learning rate of around 0.08 for 120 trials<sup>62</sup>. This may partly explain why in  
649 our case of 32 trials per condition and cue, optimal learning rates called for a relatively high optimal  
650 learning rate of 0.29, while in other studies, optimal learning rates may be lower. Regarding differences  
651 in the analysis approach, the hierarchical bayesian estimation approach used in our study produces more  
652 reliable results in comparison to maximum likelihood estimation<sup>49</sup>, which had been used in some of the

653 previous adult studies and may have led to biased results towards extreme values. Taken together, our  
654 study underscores the importance of using longitudinal data to examine developmental change as well  
655 as the importance of simulation-based optimal parameters to interpret the direction of developmental  
656 change.

657

658 Despite a relatively immature hippocampal structure in middle childhood, our results confirmed a  
659 longitudinally stable association between hippocampal volume and delayed feedback learning. However,  
660 episodic memory in this learning condition was not enhanced. This suggests a developmentally early  
661 hippocampal contribution to value-based learning during delayed feedback, which does not modulate  
662 episodic memory as much as compared to adults. Therefore, our study partially extends the findings  
663 from the adult literature to middle childhood<sup>5,12-14</sup>. The reduced effect of delayed feedback on episodic  
664 memory may be due to the protracted development of hippocampal maturation. In an aging study with  
665 a similar task, older adults failed to exhibit enhanced episodic memory for objects presented during  
666 delayed feedback trials, and they showed no enhanced hippocampal activation during delayed feedback  
667 and<sup>14</sup>. Therefore, the findings converge nicely at both childhood and older adulthood, during which the  
668 structural and functional integrity of hippocampus are known to be less optimal than at younger  
669 adulthood<sup>63-65</sup>.

670 Our brain-cognition links were only partially confirmed, as striatal volumes exhibited associations with  
671 not just immediate learning scores, as we predicted, but also with delayed learning scores. This result  
672 suggests that the striatum may be important for value-based learning in general rather than exhibiting a  
673 selective association with immediate feedback learning. This is also what we found in an explorative  
674 analysis that related the striatum to learning rate in general and further predicted longitudinal change in  
675 learning rate (Supplemental Material 5). This overall reduced brain-behavior specificity could reflect  
676 less differentiated memory systems during development, similar to findings from aging research. Here,  
677 older adults exhibited stronger striatal and hippocampal co-activation during both implicit and explicit  
678 learning, compared to more dissociable brain-behavior relationships in younger adults<sup>66</sup>. Interestingly,  
679 even in young adults, clear dissociations between memory systems such as in non-human lesion studies  
680 are uncommon, and factors like stress modulate their cooperative interaction<sup>6,10,11,67,68</sup>. Further, there are  
681 methodological differences to previous studies that could explain why striatal volumes were not  
682 uniquely associated with immediate learning in our study. For example, previous studies related reward  
683 prediction errors to striatal and hippocampal activation<sup>5,13,14</sup>, whereas we examined individual  
684 differences in brain structure and the model-derived learning scores. Future functional neuroimaging  
685 studies with children could further clarify whether children's memory systems are indeed less  
686 differentiated and explain the attenuated modulation by feedback timing. Taken together, compared to  
687 the adult literature, our results with children showed that the hippocampal structure was associated with  
688 delayed feedback learning, but did not enhance episodic memory encoding, while the striatum generally

689 supported value-based learning. These findings point towards a developmental effect of less  
690 differentiated and more cooperative memory systems in middle childhood.

691

692 Our computational modeling results revealed a separable effect of feedback timing on inverse  
693 temperature, which suggests that the memory systems modulated learning during decision-making. The  
694 reported behavioral differences in reaction time and their correlation to the inverse temperature further  
695 support the idea of a decision-related mechanism, as we found children to respond faster during delayed  
696 feedback trials and faster responding children also exhibited more value-guided choice behavior (i.e.  
697 higher inverse temperature) during delayed compared to immediate feedback. The hippocampus may  
698 contribute to a decision-related effect in the delayed feedback condition by facilitating the encoding and  
699 retrieval of learned values<sup>69</sup>. This is in contrast to previous event-related fMRI and EEG studies  
700 reporting feedback timing modulations at value update<sup>5,13,14</sup>, which may be due to at least two reasons.  
701 First, we did not include a functional brain measure to examine its differential engagement during the  
702 choice and feedback phases. Second, in such a reinforcement learning task, disentangling model  
703 parameters from the choice and feedback phases can be challenging, such as for the inverse temperature  
704 and outcome sensitivity<sup>70</sup>. Taken together, hippocampal engagement at delayed feedback may enhance  
705 outcome sensitivity as well as facilitate choice behavior through improved retrieval of action-outcome  
706 associations. A mechanism facilitating retrieval seems especially relevant in our paradigm, where  
707 multiple cues were learned and presented in a mixed order, thus creating a high memory load. To  
708 summarize, our study results suggest that feedback timing could modulate decision-making in addition  
709 to or as alternative to a mechanism at value update. However, disentangling the effects of inverse  
710 temperature and outcome sensitivity is challenging and warrants careful interpretation. Future studies  
711 might shed new light by examining neural activations at both task phases, by additionally modeling  
712 reaction times using a drift-diffusion approach, or by choosing a task design that allows independent  
713 manipulations of these phases and associated model parameters, e.g., by using different reward  
714 magnitudes during reinforcement learning, or by studying outcome sensitivity without decision-making.

715

716 One aim of developmental investigations is to identify the emergence of brain and cognition dynamics,  
717 such as the hippocampal-dependent and striatal-dependent memory systems, which have been shown to  
718 engage during reinforcement learning depending on the delay in feedback delivery. Our longitudinal  
719 study partially confirmed these brain-cognition links in middle childhood but with less specificity as  
720 previously found in adults.

721 An early existing memory system dynamic, similar to that of adults, is relevant for applying  
722 reinforcement learning principles at different timescales. In scenarios such as in the classroom, a teacher  
723 may comment on a child's behavior immediately after the action or some moments later, in par with our  
724 experimental manipulation of 1 second versus 5 seconds. Within such short range of delay in teachers'  
725 feedback, children's learning ability during the first years of schooling may function equally well and

726 depend on the striatal-dependent memory system. However, we anticipate that the reliance on the  
727 hippocampus will become even more pronounced when feedback is further delayed for longer time.  
728 Children's capacity for learning over longer timescales relies on the hippocampal-dependent memory  
729 system, which is still under development. This knowledge could help to better structure learning  
730 according to their development. Furthermore, probabilistic learning from delayed feedback may be a  
731 potential diagnostic tool to examine the hippocampal-dependent memory system during learning in  
732 children at risk. Environmental factors such as stress<sup>11</sup> and socioeconomic status<sup>39,71</sup> have been shown  
733 to affect hippocampal structure and function and may contribute to a heightened risk for  
734 psychopathology in the long term<sup>72-74</sup>. Deficits in hippocampal-dependent learning may be particularly  
735 relevant to psychopathology since dysfunctional behavior may arise from a tendency to prioritize short-  
736 term consequences over long-term ones<sup>75,76</sup> and from the maladaptive application of previously learned  
737 behavior in inappropriate contexts<sup>77</sup>. Interestingly, poor learners showed relatively less value-based  
738 learning in favor of stronger simple heuristic strategies, and excluding them modulated the hippocampal-  
739 dependent associations to learning and memory in our results. More studies are needed to further clarify  
740 the relationship between hippocampus and psychopathology during cognitive and brain development.  
741 Another key question is whether developmental trajectories observed cross-sectionally are also  
742 confirmed by longitudinal results, such as for the learning rate and inverse temperature. Our results show  
743 developmental improvements in these learning parameters in only two years. This suggests that the  
744 initial two years of schooling constitute a dynamic period for feedback-based learning, in which  
745 contingent feedback is important in shaping behavior and development.

746

747

748

749

750

751

752

## Additional Information

753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780

Funding. This study was supported by the Jacobs Foundation [grant 2014–1151] to YLS and CH. The work of YLS was also supported by the European Union (ERC-2018-StG-PIVOTAL-758898), the Deutsche Forschungsgemeinschaft (German Research Foundation, Project ID 327654276, SFB 1315, 'Mechanisms and Disturbances in Memory Consolidation: From Synapses to Systems'), and the Hessisches Ministerium für Wissenschaft und Kunst (HMWK; project 'The Adaptive Mind').

Acknowledgments. We thank the Max Planck Institute for Human Development and all members of the Jacobs study team for their vital contribution, and all participants and family members for taking part in the study.

Conflicts of interest. The authors declare no competing financial interests.

Ethics approval. This study was approved by the “Deutsche Gesellschaft für Psychologie” ethics committee (YLS\_012015).

Availability of data and code. <https://osf.io/pju65/>

Author ORCIDs.

JF: <https://orcid.org/0000-0003-0505-0798>

LZ: <https://orcid.org/0000-0002-9586-595X>

LR: <https://orcid.org/0000-0002-0144-5605>

JJM: <https://orcid.org/0000-0002-3893-8008>

JT: <https://orcid.org/0000-0001-8166-2441>

CH: <https://orcid.org/0000-0002-6580-6326>

YLS: <https://orcid.org/0000-0001-8922-7292>



## References

- 781
- 782 1. Sutton, R. S. & Barto, A. G. *Reinforcement learning: An introduction*. (MIT press, 2018).
- 783 2. Gläscher, J., Daw, N., Dayan, P. & O’Doherty, J. P. States versus Rewards: Dissociable neural  
784 prediction error signals underlying model-based and model-free reinforcement learning.  
785 *Neuron* **66**, 585 (2010).
- 786 3. Bolenz, F., Reiter, A. M. F. & Eppinger, B. Developmental Changes in Learning:  
787 Computational Mechanisms and Social Influences. *Front. Psychol.* **0**, 2048 (2017).
- 788 4. Zhang, L. & Gläscher, J. A brain network supporting social influences in human decision-  
789 making. *Sci. Adv.* **6**, 1–20 (2020).
- 790 5. Foerde, K. & Shohamy, D. Feedback Timing Modulates Brain Systems for Learning in  
791 Humans. *J. Neurosci.* **31**, 13157–13167 (2011).
- 792 6. Packard, M. G. & Goodman, J. Factors that influence the relative use of multiple memory  
793 systems. *Hippocampus* **23**, 1044–1052 (2013).
- 794 7. Goodman, J. & Packard, M. G. Memory Systems of the Basal Ganglia. *Handb. Behav.*  
795 *Neurosci.* **24**, 725–740 (2016).
- 796 8. Davidow, J. Y., Foerde, K., Galván, A. & Shohamy, D. An Upside to Reward Sensitivity: The  
797 Hippocampus Supports Enhanced Reinforcement Learning in Adolescence. *Neuron* **92**, 93–99  
798 (2016).
- 799 9. Hartley, C. A., Nussenbaum, K. & Cohen, A. O. Interactive Development of Adaptive  
800 Learning and Memory. 1–27 (2021).
- 801 10. Packard, M. G., Goodman, J. & Ressler, R. L. Emotional modulation of habit memory: neural  
802 mechanisms and implications for psychopathology. *Curr. Opin. Behav. Sci.* **20**, 25–32 (2018).
- 803 11. Schwabe, L. & Wolf, O. T. Stress and multiple memory systems: from ‘thinking’ to ‘doing’.  
804 *Trends Cogn. Sci.* **17**, 60–68 (2013).
- 805 12. Foerde, K., Race, E., Verfaellie, M. & Shohamy, D. A role for the medial temporal lobe in  
806 feedback-driven learning: Evidence from amnesia. *J. Neurosci.* **33**, 5698–5704 (2013).
- 807 13. Hölting, G. & Mecklinger, A. Feedback timing modulates interactions between feedback  
808 processing and memory encoding: Evidence from event-related potentials. *Cogn. Affect. Behav.*  
809 *Neurosci.* **2020** **202** **20**, 250–264 (2020).
- 810 14. Lighthall, N. R., Pearson, J. M., Huettel, S. A. & Cabeza, R. Feedback-Based Learning in  
811 Aging: Contributions and Trajectories of Change in Striatal and Hippocampal Systems. *J.*  
812 *Neurosci.* **38**, 8453–8462 (2018).
- 813 15. Zhang, L., Lengersdorff, L., Mikus, N., Gläscher, J. & Lamm, C. Using reinforcement learning  
814 models in social neuroscience: Frameworks, pitfalls and suggestions of best practices. *Soc.*  
815 *Cogn. Affect. Neurosci.* **15**, 695–707 (2020).
- 816 16. Nussenbaum, K. & Hartley, C. A. Reinforcement learning across development: What insights  
817 can we draw from a decade of research? *Developmental Cognitive Neuroscience* **40**, (2019).

- 818 17. Decker, J. H., Lourenco, F. S., Doll, B. B. & Hartley, C. A. Experiential reward learning  
819 outweighs instruction prior to adulthood. *Cogn. Affect. Behav. Neurosci.* **15**, 310–320 (2015).
- 820 18. Javadi, A. H., Schmidt, D. H. K. & Smolka, M. N. Differential representation of feedback and  
821 decision in adolescents and adults. *Neuropsychologia* **56**, 280–288 (2014).
- 822 19. Palminteri, S., Kilford, E. J., Coricelli, G. & Blakemore, S. J. The Computational Development  
823 of Reinforcement Learning during Adolescence. *PLoS Comput. Biol.* **12**, 1–25 (2016).
- 824 20. Master, S. L. *et al.* Distangling the systems contributing to changes in learning during  
825 adolescence. *Dev. Cogn. Neurosci.* **41**, 100732 (2020).
- 826 21. Hauser, T. U., Iannaccone, R., Walitza, S., Brandeis, D. & Brem, S. Cognitive flexibility in  
827 adolescence: Neural and behavioral mechanisms of reward prediction error processing in  
828 adaptive decision making during development. *Neuroimage* **104**, 347–354 (2015).
- 829 22. Moutoussis, M. *et al.* Change, stability, and instability in the Pavlovian guidance of behaviour  
830 from adolescence to young adulthood. *PLoS Comput. Biol.* **14**, (2018).
- 831 23. Van Den Bos, W., Cohen, M. X., Kahnt, T. & Crone, E. A. Striatum-medial prefrontal cortex  
832 connectivity predicts developmental changes in reinforcement learning. *Cereb. Cortex* **22**,  
833 1247–1255 (2012).
- 834 24. Rodriguez Buritica, J. M., Heekeren, H. R. & van den Bos, W. The computational basis of  
835 following advice in adolescents. *J. Exp. Child Psychol.* **180**, 39–54 (2019).
- 836 25. Galván, A. The Teenage Brain: Sensitivity to Rewards. *Curr. Dir. Psychol. Sci.* **22**, 88–93  
837 (2013).
- 838 26. van Duijvenvoorde, A. C. K. *et al.* A cross-sectional and longitudinal analysis of reward-  
839 related brain activation: Effects of age, pubertal stage, and reward sensitivity. *Brain Cogn.* **89**,  
840 3–14 (2014).
- 841 27. Eckstein, M. K., Wilbrecht, L. & Collins, A. G. E. What do RL Models Measure? Interpreting  
842 Model Parameters in Cognition and Neuroscience. *Curr. Opin. Behav. Sci.* **41**, 128–137 (2021).
- 843 28. Cohen, A. O., Nussenbaum, K., Dorfman, H. M., Gershman, S. J. & Hartley, C. A. The rational  
844 use of causal inference to guide reinforcement learning strengthens with age. *npj Sci. Learn.* **5**,  
845 1–9 (2020).
- 846 29. Raznahan, A. *et al.* Longitudinal four-dimensional mapping of subcortical anatomy in human  
847 development. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 1592 (2014).
- 848 30. Wierenga, L. *et al.* Typical development of basal ganglia, hippocampus, amygdala and  
849 cerebellum from age 7 to 24. *Neuroimage* **96**, 67–72 (2014).
- 850 31. Giedd, J. N. Structural Magnetic Resonance Imaging of the Adolescent Brain. *Ann. N. Y. Acad.*  
851 *Sci.* **1021**, 77–85 (2004).
- 852 32. Uematsu, A. *et al.* Developmental Trajectories of Amygdala and Hippocampus from Infancy to  
853 Early Adulthood in Healthy Individuals. *PLoS One* **7**, e46970 (2012).
- 854 33. Giedd, J. N. *et al.* Child Psychiatry Branch of the National Institute of Mental Health

- 855 Longitudinal Structural Magnetic Resonance Imaging Study of Human Brain Development.  
856 *Neuropsychopharmacology* **40**, 43 (2015).
- 857 34. Goodman, J., Marsh, R., Peterson, B. S. & Packard, M. G. Annual research review: The  
858 neurobehavioral development of multiple memory systems--implications for childhood and  
859 adolescent psychiatric disorders. *J. Child Psychol. Psychiatry.* **55**, 582–610 (2014).
- 860 35. Goddings, A. L. *et al.* The influence of puberty on subcortical brain development. *Neuroimage*  
861 **88**, 242–251 (2014).
- 862 36. Dima, D. *et al.* Subcortical volumes across the lifespan: Data from 18,605 healthy individuals  
863 aged 3–90 years. *Hum. Brain Mapp.* 1–18 (2021). doi:10.1002/hbm.25320
- 864 37. Lavenex, P. & Banta Lavenex, P. Building hippocampal circuits to learn and remember:  
865 Insights into the development of human memory. *Behavioural Brain Research* **254**, 8–21  
866 (2013).
- 867 38. Mandolesi, L., Petrosini, L., Menghini, D., Addona, F. & Vicari, S. Children’s radial arm  
868 maze performance as a function of age and sex. *Int. J. Dev. Neurosci.* **27**, 789–797 (2009).
- 869 39. Raffington, L. *et al.* Stable longitudinal associations of family income with children’s  
870 hippocampal volume and memory persist after controlling for polygenic scores of educational  
871 attainment. *Dev. Cogn. Neurosci.* **40**, 100720 (2019).
- 872 40. Raffington, L. *et al.* Effects of stress on 6- and 7-year-old children’s emotional memory differs  
873 by gender. *J. Exp. Child Psychol.* **199**, 104924 (2020).
- 874 41. Fischl, B. FreeSurfer. *Neuroimage* **62**, 774–781 (2012).
- 875 42. Phan, T. V., Smeets, D., Talcott, J. B. & Vandermosten, M. Processing of structural  
876 neuroimaging data in young children: Bridging the gap between current practice and state-of-  
877 the-art methods. *Dev. Cogn. Neurosci.* **33**, 206–223 (2018).
- 878 43. Schoemaker, D. *et al.* Hippocampus and amygdala volumes from magnetic resonance images  
879 in children: Assessing accuracy of FreeSurfer and FSL against manual segmentation.  
880 *Neuroimage* **129**, 1–14 (2016).
- 881 44. Bates, D., Mächler, M., Bolker, B. & Walker, S. Fitting Linear Mixed-Effects Models Using  
882 lme4. *J. Stat. Softw.* **67**, 1–48 (2015).
- 883 45. Brown, V. M. *et al.* Reinforcement Learning Disruptions in Individuals with Depression and  
884 Sensitivity to Symptom Change following Cognitive Behavioral Therapy. *JAMA Psychiatry*  
885 (2021). doi:10.1001/jamapsychiatry.2021.1844
- 886 46. Stan Development Team. RStan: the R interface to Stan. R package version 2.21.2. [http://mc-](http://mc-stan.org)  
887 [stan.org](http://mc-stan.org) (2021).
- 888 47. R Core Team. R: A Language and Environment for Statistical Computing. (2021).
- 889 48. Ahn, W.-Y., Haines, N. & Zhang, L. Revealing Neurocomputational Mechanisms of  
890 Reinforcement Learning and Decision-Making With the hBayesDM Package. *Comput.*  
891 *Psychiatry* **1**, 24 (2017).

- 892 49. Brown, V. M., Chen, J., Gillan, C. M. & Price, R. B. Improving the Reliability of  
893 Computational Analyses: Model-Based Planning and Its Relationship With Compulsivity. *Biol.*  
894 *Psychiatry Cogn. Neurosci. Neuroimaging* **5**, 601–609 (2020).
- 895 50. Vehtari, A., Gelman, A. & Gabry, J. Practical Bayesian model evaluation using leave-one-out  
896 cross-validation and WAIC. *Stat. Comput.* **27**, 1413–1432 (2017).
- 897 51. Yao, Y., Vehtari, A., Simpson, D. & Gelman, A. Using Stacking to Average Bayesian  
898 Predictive Distributions (with Discussion). *Bayesian Anal.* **13**, 917–1007 (2018).
- 899 52. Crawley, D. *et al.* Modeling flexible behavior in childhood to adulthood shows age-dependent  
900 learning mechanisms and less optimal learning in autism in each age group. *PLoS Biol.* **18**, 1–  
901 25 (2020).
- 902 53. Wilson, R. C. & Collins, A. G. E. Ten simple rules for the computational modeling of  
903 behavioral data. *Elife* **8**, 1–33 (2019).
- 904 54. Kievit, R. A. *et al.* Developmental cognitive neuroscience using latent change score models: A  
905 tutorial and applications. *Dev. Cogn. Neurosci.* **33**, 99–117 (2018).
- 906 55. Ferrer, E. & McArdle, J. J. Longitudinal modeling of developmental changes in psychological  
907 research. *Curr. Dir. Psychol. Sci.* **19**, 149–154 (2010).
- 908 56. Sluis, S. van der, Verhage, M., Posthuma, D. & Dolan, C. V. Phenotypic Complexity,  
909 Measurement Bias, and Poor Phenotypic Resolution Contribute to the Missing Heritability  
910 Problem in Genetic Association Studies. *PLoS One* **5**, e13929 (2010).
- 911 57. Little, T. *Longitudinal structural equation modeling*. (Guilford Press, 2013).  
912 doi:10.2/JQUERY.MIN.JS
- 913 58. Chierchia, G. *et al.* Confirmatory reinforcement learning changes with age during adolescence.  
914 *Dev. Sci.* e13330 (2021). doi:10.1111/desc.13330
- 915 59. Habicht, J., Bowler, A., Moses-Payne, M. E. & Hauser, T. U. Children are full of optimism, but  
916 those rose-tinted glasses are fading – reduced learning from negative outcomes drives  
917 hyperoptimism in children. (2021).
- 918 60. Overman, W. H. Sex differences in early childhood, adolescence, and adulthood on cognitive  
919 tasks that rely on orbital prefrontal cortex. *Brain Cogn.* **55**, 134–147 (2004).
- 920 61. Evans, K. L. & Hampson, E. Sex-dependent effects on tasks assessing reinforcement learning  
921 and interference inhibition. *Front. Psychol.* **6**, 1–10 (2015).
- 922 62. Behrens, T. E. J., Woolrich, M. W., Walton, M. E. & Rushworth, M. F. S. Learning the value  
923 of information in an uncertain world. *Nat. Neurosci.* **10**, 1214–1221 (2007).
- 924 63. Shing, Y. L. *et al.* Episodic memory across the lifespan: The contributions of associative and  
925 strategic components. *Neurosci. Biobehav. Rev.* **34**, 1080–1091 (2010).
- 926 64. Keresztes, A. *et al.* Hippocampal maturity promotes memory distinctiveness in childhood and  
927 adolescence. *Proc. Natl. Acad. Sci. U. S. A.* **114**, 9212–9217 (2017).
- 928 65. Ghetti, S. & Bunge, S. A. Neural changes underlying the development of episodic memory

- 929 during middle childhood. *Dev. Cogn. Neurosci.* **2**, 381–395 (2012).
- 930 66. Dennis, N. A. & Cabeza, R. Age-related dedifferentiation of learning systems: An fMRI study  
931 of implicit and explicit learning. *Neurobiol. Aging* **32**, 2318.e17–2318.e30 (2011).
- 932 67. Ferbinteanu, J. Contributions of Hippocampus and Striatum to Memory-Guided Behavior  
933 Depend on Past Experience. *J. Neurosci.* **36**, 6459–6470 (2016).
- 934 68. White, N. M. & McDonald, R. J. Multiple Parallel Memory Systems in the Brain of the Rat.  
935 *Neurobiol. Learn. Mem.* **77**, 125–184 (2002).
- 936 69. Shadlen, M. N. N. & Shohamy, D. Decision Making and Sequential Sampling from Memory.  
937 *Neuron* **90**, 927–939 (2016).
- 938 70. Browning, M., Paulus, M. & Huys, Q. J. M. What is computational psychiatry good for? *Biol.*  
939 *Psychiatry* **0**, (2022).
- 940 71. Hackman, D. A., Farah, M. J. & Meaney, M. J. Socioeconomic status and the brain:  
941 mechanistic insights from human and animal research. *Nat. Rev. Neurosci.* **2010 119 11**, 651–  
942 659 (2010).
- 943 72. Frodl, T. *et al.* Childhood stress, serotonin transporter Gene and Brain structures in major  
944 depression. *Neuropsychopharmacology* **35**, 1383–1390 (2010).
- 945 73. Lucassen, P. J., Korosi, A., Krugers, H. J. & Oomen, C. A. *Early Life Stress- and Sex-*  
946 *Dependent Effects on Hippocampal Neurogenesis. Stress: Neuroendocrinology and*  
947 *Neurobiology* **2**, (Elsevier Inc., 2017).
- 948 74. Rahman, M. M., Callaghan, C. K., Kerskens, C. M., Chattarji, S. & O’Mara, S. M. Early  
949 hippocampal volume loss as a marker of eventual memory deficits caused by repeated stress.  
950 *Sci. Rep.* **6**, 1–15 (2016).
- 951 75. Levin, M. E., Haeger, J., Ong, C. W. & Twohig, M. P. An Examination of the Transdiagnostic  
952 Role of Delay Discounting in Psychological Inflexibility and Mental Health Problems. *Psychol.*  
953 *Rec.* **68**, 201–210 (2018).
- 954 76. Von Siebenthal, Z. *et al.* Decision-making impairments following insular and medial temporal  
955 lobe resection for drug-resistant epilepsy. *Soc. Cogn. Affect. Neurosci.* **12**, 128–137 (2017).
- 956 77. Maren, S., Phan, K. L. & Liberzon, I. The contextual brain: Implications for fear conditioning,  
957 extinction and psychopathology. *Nat. Rev. Neurosci.* **14**, 417–428 (2013).

958