

1 **A new model of decision processing in instrumental learning tasks**

2

3 Steven Miletić^{a*}, Russell J. Boag^a, Anne C. Trutti^{a,b}, Birte U. Forstmann^a, Andrew Heathcote^{a,c}

4

5

6 ^a University of Amsterdam, Department of Psychology, Nieuwe Achtergracht 129B, Amsterdam, The

7 Netherlands

8 ^b Leiden University, Department of Psychology, Wassenaarseweg 52, Leiden, The Netherlands

9 ^c University of Newcastle, School of Psychology, Newcastle, Australia

10

11 *Correspondence concerning this article should be addressed to Steven Miletić, Nieuwe Achtergracht 129B,

12 1001NK Amsterdam, The Netherlands. E: s.miletic@uva.nl

13

14 **Abstract**

15 Learning and decision making are interactive processes, yet cognitive modelling of error-
16 driven learning and decision making have largely evolved separately. Recently, evidence
17 accumulation models (EAMs) of decision making and reinforcement learning (RL) models of
18 error-driven learning have been combined into joint RL-EAMs that can in principle address
19 these interactions. However, we show that the most commonly used combination, based on the
20 diffusion decision model (DDM) for binary choice, consistently fails to capture crucial aspects
21 of response times observed during reinforcement learning. We propose a new RL-EAM based
22 on an advantage racing diffusion (ARD) framework for choices among two or more options
23 that not only addresses this problem but captures stimulus difficulty, speed-accuracy trade-off,
24 and stimulus-response-mapping reversal effects. The RL-ARD avoids fundamental limitations
25 imposed by the DDM on addressing effects of absolute values of choices, as well as extensions
26 beyond binary choice, and provides a computationally tractable basis for wider applications.

27

28 *Keywords:* Decision making, reinforcement learning, evidence-accumulation models, speed-
29 accuracy trade-off, reversal learning.

30

31

32

33

34 Learning and decision-making are mutually influential cognitive processes. Learning processes
35 refine the internal preferences and representations that inform decisions, and the outcomes of
36 decisions underpin feedback-driven learning (Bogacz and Larsen, 2011). Although this relation
37 between learning and decision-making has been acknowledged (Bogacz and Larsen, 2011;
38 Dayan and Daw, 2008), the study of cognitive processes underlying feedback-driven learning
39 on the one hand, and of perceptual and value-based decision-making on the other, have
40 progressed as largely separate scientific fields. In the study of error-driven learning (O’Doherty
41 et al., 2017; Sutton and Barto, 2018), the decision process is typically simplified to soft-max,
42 a descriptive model that offers no process-level understanding of how decisions arise from
43 representations, and ignores choice response times (RTs). In the study of decision-making
44 using evidence-accumulation models (EAMs; Donkin and Brown, 2018; Forstmann et al.,
45 2016; Ratcliff et al., 2016), tasks are typically designed to minimize the influence of learning,
46 and residual variability caused by learning is treated as noise.

47 Recent advances (Fontanesi et al., 2019a, 2019b; Luzzardo et al., 2017; McDougle and
48 Collins, 2020; Miletić et al., 2020; Millner et al., 2018; Pedersen et al., 2017; Pedersen and
49 Frank, 2020; Sewell et al., 2019; Sewell and Stallman, 2020; Shahar et al., 2019; Turner, 2019)
50 have emphasized how both modelling traditions can be combined in joint models of
51 reinforcement learning (RL) and evidence-accumulation decision-making processes, providing
52 mutual benefits for both fields. Combined models generally propose that value-based decision-
53 making and learning interact as follows: For each decision a subject gradually accumulates
54 evidence for each choice option by sampling from a distribution of memory representations of
55 the subjective value (or *expected reward*) associated with each choice option (known as *Q-*
56 *values*). Once a threshold level of evidence is reached, they commit to the decision and initiate
57 a corresponding motor process. The response triggers feedback, which is used to update the
58 internal representation of subjective values. The next time the subject encounters the same
59 choice options, this updated internal representation changes evidence accumulation.

60 The RL-EAM framework has many benefits (Miletić et al., 2020). It allows for studying a
61 rich set of behavioral data simultaneously, including entire RT distributions and trial-by-trial
62 dependencies in choices and RTs. It posits a theory of evidence accumulation that assumes a
63 memory representation of rewards is the source of evidence, and it formalizes how these
64 memory representations change due to learning. It complements earlier work connecting
65 theories of reinforcement learning and decision-making (Bogacz and Larsen, 2011; Dayan and
66 Daw, 2008) and their potential neural implementation in basal ganglia circuits (Bogacz and
67 Larsen, 2011), by presenting a measurement model that can be fit to, and makes predictions
68 about, behavioral data. Adding to benefits in terms of theory building, the RL-EAM framework
69 also has potential to improve parameter recovery properties compared to standard RL models
70 (Shahar et al., 2019), and allows for the estimation of single-trial parameters of the decision
71 model, which can be crucial in the analysis of neuroimaging data.

72 An important challenge of this framework is the number of modeling options in both the
73 fields of reinforcement learning and decision-making. Even considering only model-free (as
74 opposed to model-based (Daw and Dayan, 2014)) reinforcement learning, there exists a variety
75 of learning rules (e.g., Palminteri et al., 2015; Rescorla and Wagner, 1972; Rummery and
76 Niranjan, 1994; Sutton, Richard, 1988), as well as the possibility of multiple learning rates for
77 positive and negative prediction errors (Christakou et al., 2013; Daw et al., 2002; Frank et al.,

78 2009; Gershman, 2015; Haughey et al., 2007; Niv et al., 2012), and many additional concepts,
79 such as eligibility traces to allow for updating of previously visited states (Barto et al., 1981;
80 Bogacz et al., 2007). Similarly, in the decision-making literature, there exists a wide range of
81 evidence-accumulation models, including most prominently the diffusion decision model
82 (DDM; Ratcliff, 1978; Ratcliff et al., 2016) and race models such as the linear ballistic
83 accumulator model (LBA; Brown and Heathcote, 2008) and racing diffusion (RD) models
84 (Boucher et al., 2007; Hawkins and Heathcote, 2020; Leite and Ratcliff, 2010; Logan et al.,
85 2014; Purcell et al., 2010; Ratcliff et al., 2011; Tillman et al., 2020).

86 The existence of this wide variety of modelling options is a double-edged sword. On the one
87 hand, it highlights the success of the general principles underlying both modelling traditions
88 (i.e., learning from prediction errors and accumulate-to-threshold decisions) in explaining
89 behavior, and it allows for studying specific learning/decision-making phenomena. On the
90 other hand, it constitutes a bewildering combinatorial explosion of potential RL-EAMs; here
91 we provide empirical grounds to navigate this problem with respect to EAMs.

92 The DDM is the dominant EAM as currently used in reinforcement learning (Fontanesi et
93 al., 2019a, 2019b; Millner et al., 2018; Pedersen et al., 2017; Pedersen and Frank, 2020; Sewell
94 et al., 2019; Sewell and Stallman, 2020; Shahar et al., 2019), but this choice is without
95 experimental justification. Furthermore, the DDM has several theoretical drawbacks, such as
96 its inability to explain multi-alternative decision-making and its strong commitment to the
97 accumulation of the evidence *difference*, which leads to difficulties in explaining behavioral
98 effects of absolute stimulus and reward magnitudes without additional mechanisms (Fontanesi
99 et al., 2019a; Ratcliff et al., 2018; Teodorescu et al., 2016). Here, we compare the performance
100 of different decision-making models in explaining choice behavior in a variety of instrumental
101 learning tasks. Models that fail to capture crucial aspects of performance run the risk of
102 producing misleading psychological inferences. For EAMs, the full RT distribution (i.e., its
103 level of variability and skew) have proven to be crucial. Hence, it is important to assess which
104 RL-EAMs are able to capture not only learning-related changes in choice probabilities and
105 mean RT, but also the general shape of the entire RT distribution and how it changes with
106 learning. Further, in order to be held forth as a general modeling framework, it is important to
107 capture how all of these measures interact with key phenomena in the decision-making and
108 learning literature.

109 We compare the RL-DDM with two RL-EAMs based on a racing accumulator architecture
110 (Figure 1). All RL-EAMs assume evidence accumulation is driven by Q-values, which change
111 based on error-driven learning as governed by the classical State-Action-Reward-State-Action
112 (SARSA; Rummery and Niranjana, 1994) update rule. Rather than a two-sided DDM process
113 (Figure 1A), the alternative models adopt a neurally plausible RD architecture (Ratcliff et al.,
114 2007), which conceptualize decision making as a statistically independent race between single-
115 sided diffusive accumulators, each collecting evidence for a different choice option. The first
116 accumulator to reach its threshold triggers motor processes that execute the corresponding
117 decision. The alternative models differ in how the mean values of evidence are constituted. The
118 first model, the RL-RD (Figure 1B), postulates accumulators are driven by the expected reward
119 for their choice, plus a stimulus-independent baseline (c.f. an *urgency* signal; Miletic and Van
120 Maanen, 2019). The second model, the RL-ARD (advantage racing diffusion), uses the recently
121 proposed *advantage* framework (Van Ravenzwaaij et al., 2020), assuming that each

122 accumulator is driven by weighted combination of three terms: the *difference* (“advantage”) in
 123 mean reward expectancy of one choice option over the other, the *sum* of the mean reward
 124 expectancies, and the urgency signal. In perceptual choice the advantage term consistently
 125 dominates the sum term by an order of magnitude (Van Ravenzwaaij et al., 2020), but the sum
 126 term is necessary to explain the effects of absolute stimulus magnitude. We also fit a limited
 127 version of this model, RL-IARD, with the weight of the sum term set to zero to test whether
 128 accounting for the influence of the sum is necessary even when reward magnitude is not
 129 manipulated, as was the case in our experiments. The importance of sum and advantage terms
 130 is also quantified by their weights as estimated in full RL-ARD model fits.

131 For all models, we first test how well they account for RT distributions (central tendency,
 132 variability, and skewness of RTs), accuracies, and learning-related changes in RT distributions
 133 and accuracies in a typical instrumental learning task (Frank, 2004). In this experiment we also
 134 manipulated difficulty, that is, the magnitude of the difference in average reward between pairs
 135 of options. In two further experiments we test the ability of the RL-EAMs to capture key
 136 behavioral phenomena in the decision-making and reinforcement-learning literatures,
 137 respectively, speed-accuracy trade-off (SAT), and reversals in reward contingencies. Again,
 138 these tests required a comprehensive account of not only choice probabilities but also the full
 139 distribution of RT, and learning-related changes thereof.

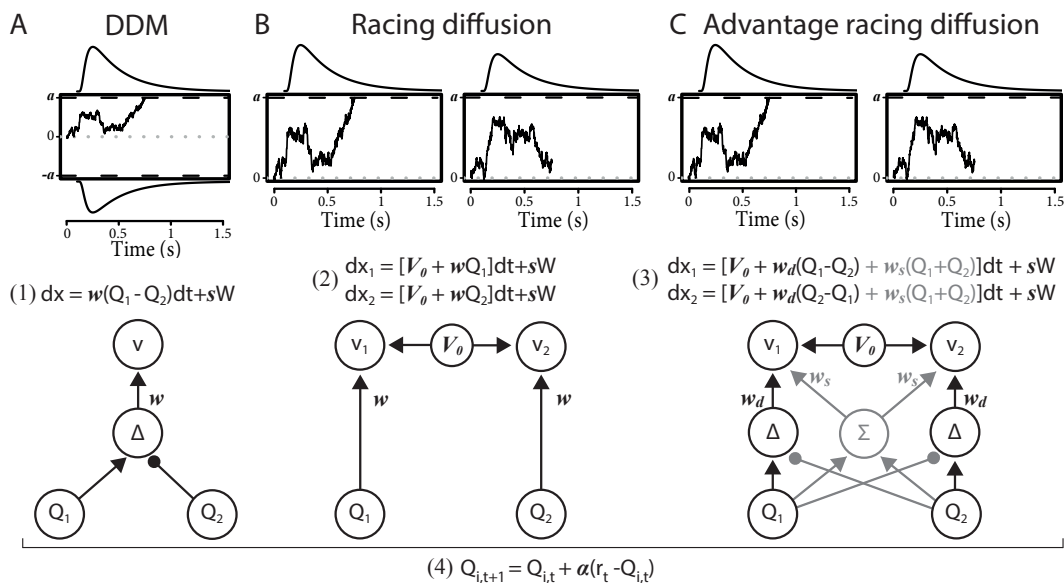


Figure 1. Comparison of the decision-making models. Bottom graphs visualize how Q-values are linked to accumulation rates. Top panel illustrates the evidence-accumulation process of the DDM (panel A) and racing diffusion (RD) models (panels B and C). Note that in the race models there is no lower bound. Equations 1-3 formally link Q-values to evidence-accumulation rates. In the RL-DDM, the difference (Δ) in Q-values is accumulated, weighted by free parameter w , plus additive within-trial white noise W with standard deviation s . In the RL-RD, the (weighted) Q-values for both choice options are independently accumulated. An evidence-independent baseline urgency term, V_0 (equal for all accumulators), further drives evidence accumulation. In the RL-ARD models, the advantages (Δ) in Q-values are accumulated as well, plus the evidence-independent baseline term V_0 . The grey icons indicate the influence of the Q-value *sum* (Σ) on evidence accumulation, which is not included in the limited variant of the RL-ARD. In all panels, bold-italic faced characters indicate parameters. Q_1 and Q_2 are Q-values for both choice options, which are updated according to a SARSA learning rule (equation (4) at the bottom of the graph), with learning rate α .

140

141 **Results**

142 In the first experiment, participants made decisions between four sets of two abstract choice
143 stimuli, each associated with a fixed reward probability (Figure 2A). On each trial, one choice
144 option always had a higher expected reward than the other; we refer to this choice as the
145 ‘correct’ choice. After each choice, participants received feedback in the form of points.
146 Reward probabilities, and therefore choice difficulty, differed between the four sets (Figure
147 2B). In total, data from 55 subjects were included in the analysis, each performing 208 trials
148 (see methods).

149 Throughout, we summarize RT distributions by calculating the 10th, 50th (median) and 90th
150 percentiles separately for correct and error responses. The median summarizes central
151 tendency, the difference between 10th and 90th percentiles summarizes variability and the larger
152 difference between the 90th and 50th percentiles than between the 50th and 10th percentiles
153 summarizes the positive skew that is always observed in RT distributions. To visualize the
154 effect of learning, we divided all trials in 10 bins (approximately 20 trials each), and calculated
155 accuracy and the RT percentiles per bin. Note that model fitting was not based on these data
156 summaries. Instead, we used hierarchical Bayesian methods to fit models to the data from every
157 trial and participant simultaneously. We compared model fits informally using posterior
158 predictive distributions—calculating the same summary statistics on data generated from the
159 fitted model as we did for the empirical data—and formally using the Bayesian Predictive
160 Information Criterion (BPIC; Ando, 2007). The former method allows us to assess the absolute
161 quality of fit (Palminteri et al., 2017) and detect misfits; the latter provides a model-selection
162 criterion that trades off quality of fit with model complexity (lower BPICs are preferred),
163 ensuring that a better fit is not only due to greater model flexibility.
164

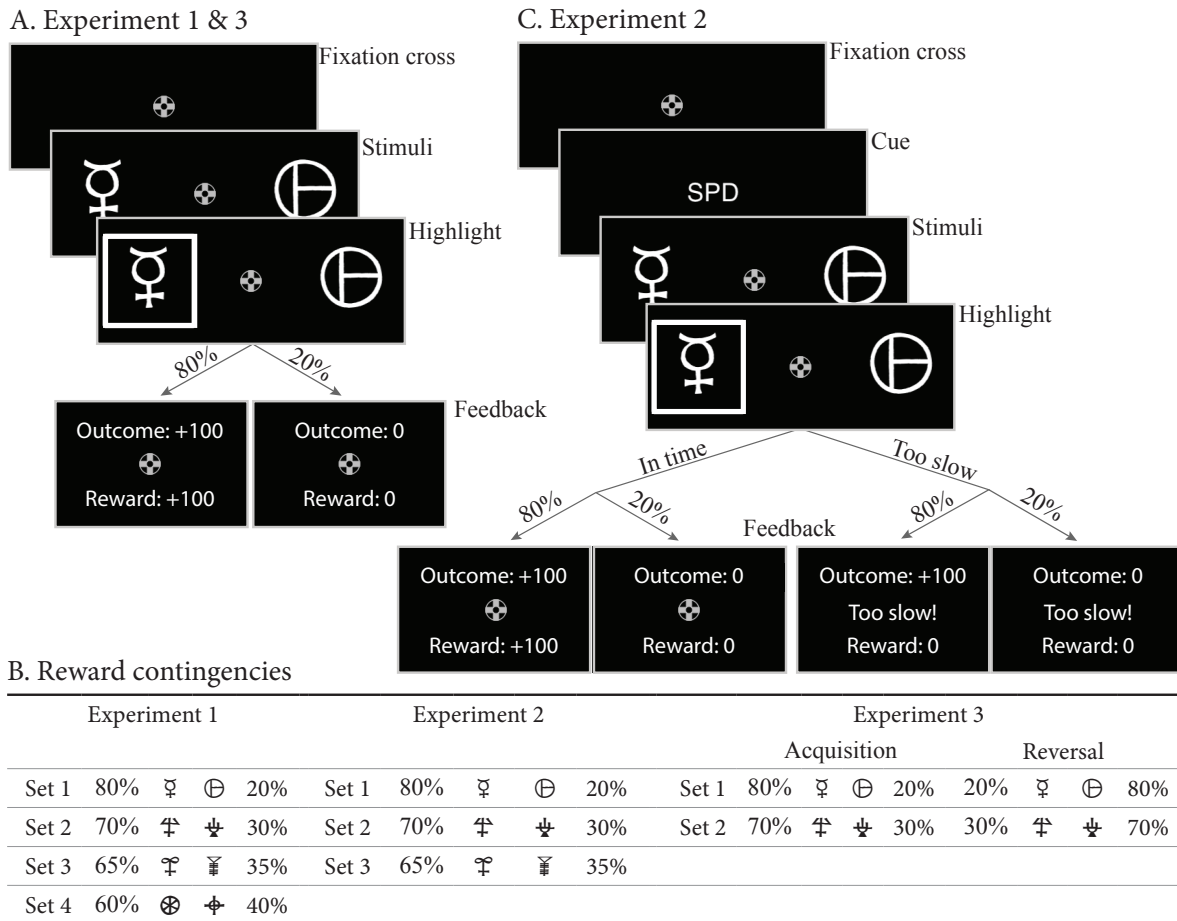


Figure 2. Paradigms for all experiments. A: Example trial for experiment 1 and 3. Each trial starts with a fixation cross, followed by the presentation of the stimulus (until choice is made or 2.5 s elapses), a brief highlight of the chosen option, and probabilistic feedback. Reward probabilities are summarized in B. Percentages indicate the probabilities of receiving +100 points for a choice (with 0 otherwise). The actual symbols used differed between experiments and participants. In experiment 3, the acquisition phase lasted 61-68 trials (uniformly sampled each block), after which the reward contingencies for each stimulus set reversed. C: Example trial for experiment 2, which added a cue prior to each trial ('SPD' or 'ACC'), and had feedback contingent on both the choice and choice timing. In the SPD condition, RTs under 700 ms were considered in time, and too slow otherwise. In the ACC condition, choices were in time as long as they were made in the stimulus window of 1.5 s. Positive feedback "Outcome: +100" and "Reward: +100" were shown in green letters, negative feedback ("Outcome: 0", "Reward: 0", and "Too slow!") were shown in red letters.

165

166 We first examine results aggregated over difficulty conditions. The posterior predictives of all
 167 four RL-EAMs are shown in Figure 3, with the top row showing accuracies, and the middle
 168 and bottom rows correct and error RT distributions (parameter estimates for all models can be
 169 found in the Table 1). The RL-DDM generally explains the learning-related increase in
 170 accuracy well, and if only the central tendency were relevant it might be considered to provide
 171 an adequate account of RT, although correct median RT is systematically under-estimated.
 172 However, RT variability and skew are severely over-estimated. The RL-RD largely overcomes

173 the RT distribution misfit, but it overestimates RTs in the first trial bins, and while capturing
174 an increase in accuracy over trials, it is systematically underestimated. The RL-ARD models
175 provide the best explanation of all key aspects of the data: except for a slight underestimation
176 of accuracy in early trial bins (largely shared with the RL-DDM), they capture accuracy well,
177 and like the RL-RD, they capture the RT distributions well, but without overpredicting the RTs
178 in the early trials. The two RL-ARD models do not differ greatly in fit, except that the limited
179 version slightly underestimates the decrease in RT with learning.
180

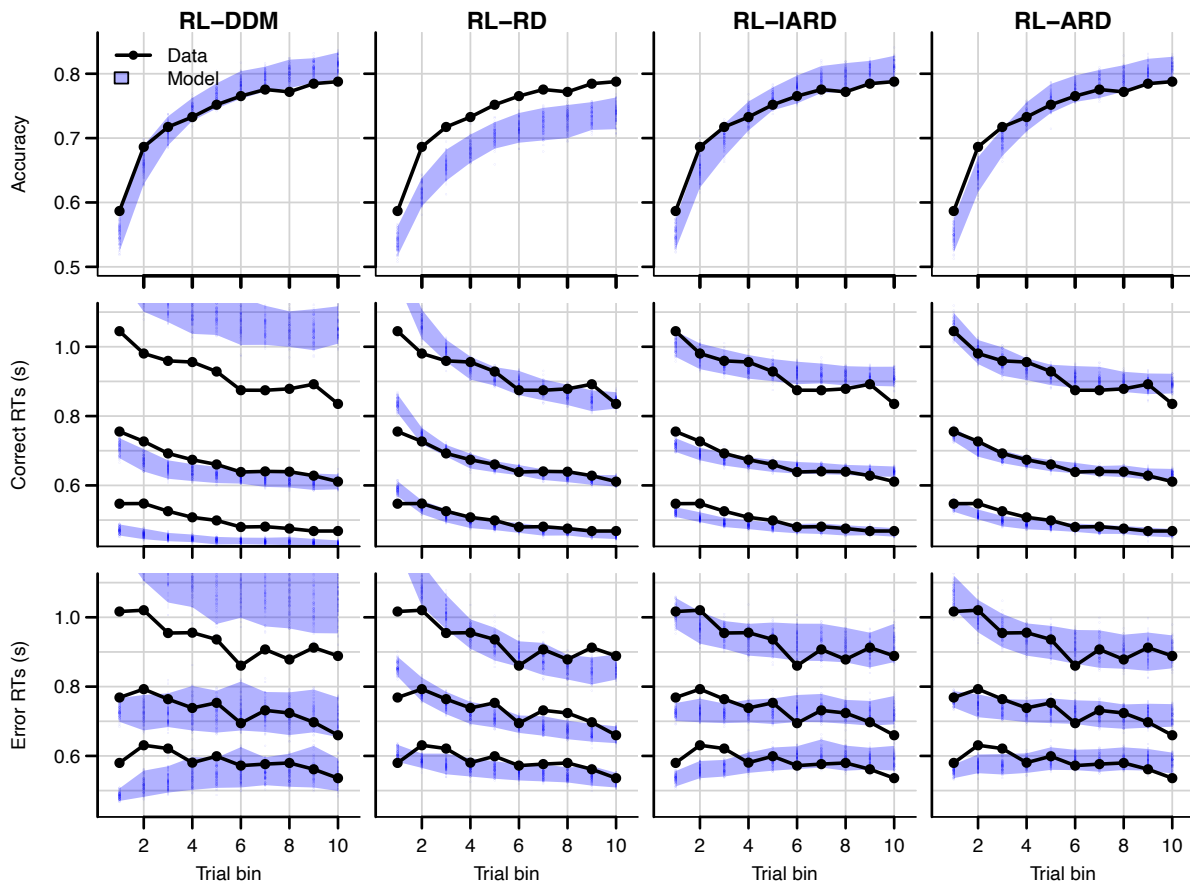


Figure 3. Comparison of posterior predictive distributions of the four RL-EAMs. Data (black) and posterior predictive distribution (blue) of the RL-DDM (left column), RL-RD, RL-IARD, and RL-ARD (right column). Top row depicts accuracy over trial bins. Middle and bottom row show 10th, 50th, and 90th RT percentiles for the correct (middle row) and error (bottom row) response over trial bins. Shaded areas correspond to the 95% credible interval of the posterior predictive distributions. All data are collapsed across participants and difficulty conditions.

181
182
183

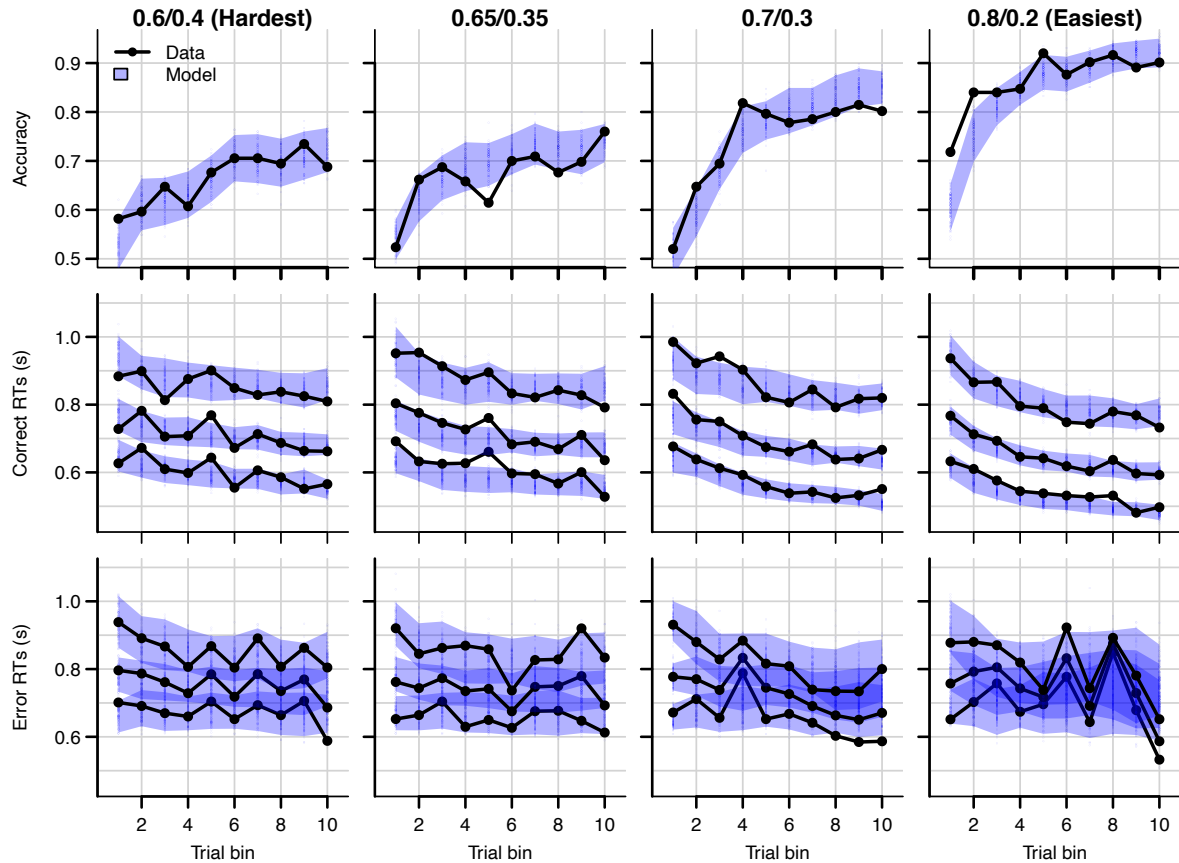


Figure 4. Data (black) and posterior predictive distribution of the RL-ARD (blue), separately for each difficulty condition. Column titles indicate the reward probabilities, with 0.6/0.4 being the most difficult, and 0.8/0.2 the easiest condition. Top row depicts accuracy over trial bins. Middle and bottom rows show 10th, 50th, and 90th RT percentiles for the correct (middle row) and error (bottom row) response over trial bins. Shaded areas correspond to the 95% credible interval of the posterior predictive distributions. All data and fits are collapsed across participants.

184

185 Figure 4 shows the data and RL-ARD model fit separated by difficulty (see Figure 4-figure
 186 supplement 1 for equivalent RL-DDM fits, which again fail to capture RT distributions). The
 187 RL-ARD model displays the same excellent fit as to data aggregated over difficulty, except
 188 that it underestimates accuracy in early trials in the easiest condition (Figure 4, bottom right
 189 panel). Further inspections of the data revealed that 17 participants (31%) reached perfect
 190 accuracy in the first bin in this condition. Likely, they guessed correctly on the first occurrence
 191 of the easiest choice pair, repeated their choice, and received too little negative feedback in the
 192 next repetitions to change their choice strategy. Figure 4-figure supplement 2 shows that, with
 193 these 17 participants removed, the overestimation is largely mitigated. SARSA assumes
 194 learning from feedback, and so cannot explain such high early accuracies. Working memory
 195 processes could have aided performance in the easiest condition, since the total number of
 196 stimuli pairs was limited and feedback was quite reliable, making it relatively easy to remember
 197 correct-choice options (Collins and Frank, 2018, 2012a; McDougle and Collins, 2020).

198

199 **Reward magnitude and Q-value evolution**

200 Q-values represent the participants' internal beliefs about how rewarding each choice option
201 is. The RL-IARD and RL-DDM assume drift rates are driven only by the difference in Q-values
202 (Figure 5), and both underestimate the learning-related decrease in RTs. Similar RL-DDM
203 underestimation has been detected before (Pedersen et al., 2017), with the proposed remedy
204 being a decrease in the decision bound with time (but with no account of RT distributions).
205 The RL-ARD explains the additional speed-up through the increasing *sum* of Q-values over
206 trials (Figure 5C), which in turn increases drift rates (Figure 5D). In line with observations in
207 perceptual decision-making (Van Ravenzwaaij et al., 2020), the effect of the expected reward
208 magnitude on drift rate is smaller (on average, $w_s = 0.36$) than that of the Q-value difference
209 ($w_D = 2.25$) and the urgency signal ($V_0 = 2.45$). Earlier work using an RL-DDM (Fontanesi
210 et al., 2019a) showed that higher reward magnitudes decrease RTs in reinforcement learning
211 paradigms. There, the reward magnitude effect on RT was accounted for by allowing the
212 threshold to change as a function of magnitude. However, this requires participants to rapidly
213 adjust their threshold based on the identity of the stimuli, something that is usually not
214 considered possible in EAMs (Donkin et al., 2011; Ratcliff, 1978). The RL-ARD avoids this
215 problem, with magnitude effects entirely mediated by drift rates, and our result show expected
216 reward magnitudes influence RTs due to learning even in the absence of a reward magnitude
217 manipulation. Because the sum affects each accumulator equally, it changes RT with little
218 effect on accuracy.

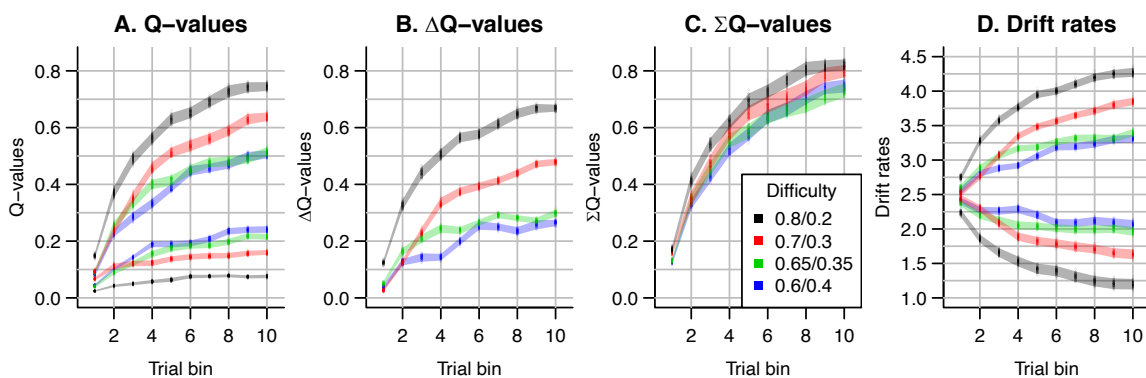


Figure 5. The evolution of Q-values and their effect on drift rates in the RL-ARD. A depicts raw Q-values, separate for each difficulty condition (colors). B and C depict the Q-value differences and the Q-value sums over time. The drift rates (D) are a weighted sum of the Q-value differences and Q-value sums, plus an intercept.

219

220 Speed-accuracy trade-off

221 Speed-accuracy trade-off (SAT) refers to the ability to strategically trade-off decision speed
222 for decision accuracy (Bogacz et al., 2010; Pachella and Pew, 1968; Ratcliff and Rouder,
223 1998). As participants can voluntarily trade speed for accuracy, RT and accuracy are not
224 independent variables, so analysis methods considering only one of these variables while
225 ignoring the other (e.g., soft-max, which only focuses on choice accuracy) can be misleading.
226 EAMs simultaneously consider RTs and accuracy and allow for estimation of SAT settings.
227 The classical explanation in the DDM framework (Ratcliff and Rouder, 1998) holds that
228 participants adjust their SAT by changing the decision threshold: increasing thresholds require
229 a participant to accumulate more evidence, leading to slower but more accurate responses.

230 Empirical work draws a more complex picture. Several papers suggest that in addition to
231 thresholds, drift rates (Arnold et al., 2015; Heathcote and Love, 2012; Ho et al., 2012; Rae et
232 al., 2014; Sewell and Stallman, 2020) and sometimes even non-decision times (Arnold et al.,
233 2015; Voss et al., 2004) can be affected. Increases in drift rates in a race model could indicate
234 an urgency signal, implemented by drift gain modulation, with qualitatively similar effects to
235 collapsing thresholds over the course of a decision (Cisek et al., 2009; Hawkins et al., 2015;
236 Miletic, 2016; Miletic and Van Maanen, 2019; Murphy et al., 2016; Thura and Cisek, 2016;
237 Trueblood et al., 2020; van Maanen et al., 2019). In cognitively demanding tasks, it has been
238 shown that two distinct components of evidence accumulation (quality and quantity of
239 evidence) are affected by SAT manipulations, with quantity of evidence being analogous to an
240 urgency signal (Boag et al., 2019b, 2019a). Recent evidence suggests that different SAT
241 manipulations can affect different psychological processes: cue-based manipulations that
242 instruct participants to be fast or accurate, lead to overall threshold adjustments, whereas
243 deadline-based manipulations lead to a collapse of thresholds (Katsimpokis et al., 2020).

244 Here, we apply an SAT manipulation in an instrumental learning task (Figure 2C). This
245 paradigm differed from experiment 1 by the inclusion of a cue-based instruction to either stress
246 response *speed* ('SPD') or response *accuracy* ('ACC') prior to each choice (randomly
247 interleaved). Furthermore, on speed trials, participants had to respond within 0.7 s to receive a
248 reward. Feedback was determined based on both the choice's probabilistic outcome ('+100' or
249 '+0') and the RT: On trials where participants responded too late, they were additionally
250 informed of the reward associated with their choice, had they been in time, so that they always
251 received the feedback required to learn from their choices. After exclusions (see methods), data
252 from 19 participants (324 trials each) were included in the analyses.

253 We used two mixed effects models to confirm the effect of the manipulation. A linear model
254 predicting RT confirmed an interaction between trial bin and cue ($b = 1.954 \cdot 10^{-4}$, $SE = 2.20 \cdot 10^{-3}$,
255 95% CI [$1.52 \cdot 10^{-2}$, $2.38 \cdot 10^{-2}$], $p < 10^{-16}$), a main effect of cue ($b = -1.913 \cdot 10^{-1}$, $SE = 9.81 \cdot 10^{-3}$,
256 95% CI [-0.21, -0.17], $p < 10^{-16}$) and a main effect of trial bin ($b = -2.21 \cdot 10^{-4}$, $SE = 1.55 \cdot 10^{-3}$,
257 95% CI [$-2.51 \cdot 10^{-3}$, $-1.9 \cdot 10^{-3}$], $p < 10^{-16}$). Thus, RTs decreased with trial bin, were faster for
258 the speed condition, but the effect of the cue was smaller for later trial bins. The logistic mixed
259 effects model predicting choice accuracy showed a main effect of the cue ($\beta = -0.37$, $SE =$
260 0.12 , std. $\beta = -0.22$, $p < 0.01$) and trial bin ($\beta = 0.45$, $SE = 0.07$, std. $\beta = 0.29$, $p <$
261 0.001), but not for an interaction ($\beta = 0.12$, $SE = 0.09$, std. $\beta = 0.08$, $p = 0.174$). Hence,
262 participants were more accurate in the accuracy condition, and there was an increase in
263 accuracy over trial bins, but there was no evidence for a difference in the increase over trial
264 bins between SAT conditions.

265 To illustrate the importance of simultaneously analyzing RTs and choice behavior, we first
266 test whether a soft-max model (which ignores RTs) is able to capture the behavioral changes
267 in choice probability due to the manipulation. We fit two soft-max models to the data: One
268 with a single inverse temperature parameter, and one with an inverse temperature parameter
269 per SAT condition. The soft-max model with separate parameters per condition was
270 outperformed by a model with a single parameter ($\Delta BPIC = 11$), indicating that a researcher
271 using soft-max would have concluded that there was no difference in choice behavior between
272 conditions. Clearly, the difference in accuracy (and RTs) did indicate there were differences in

273 behavior (as statistically confirmed above), showing that soft-max fails to capture a strong and
274 well-known phenomenon of decision-making.

275 Next, we compared the RL-DDM and RL-ARD, and in light of the multiple psychological
276 mechanisms potentially affected by the SAT manipulation, we allowed different combinations
277 of threshold, drift rate, and for the RL-ARD urgency, to vary with the SAT manipulation. We
278 fit three RL-DDM models, varying either threshold, the Q-value weighting on the drift rates
279 parameter (Sewell and Stallman, 2020), or both. For the RL-ARD, we fit all seven possible
280 models with different combinations of the threshold, urgency, and drift rate parameters free to
281 vary between SAT conditions.

282 Formal model comparison (see Table 1 for all BPIC values) indicated that the RL-ARD
283 model combining response caution and urgency effects provides the best explanation of the
284 data, in line with earlier research in non-learning contexts (Katsimpokis et al., 2020; Miletic
285 and Van Maanen, 2019; Rae et al., 2014; Thura and Cisek, 2016). The advantage for the RL-
286 ARD was substantial; the best RL-DDM (with only a threshold effect) performed worse than
287 the worst RL-ARD model. The data and posterior predictive distributions of the best RL-DDM
288 model and the winning RL-ARD model are shown in Figure 6. As in experiment 1, the RL-
289 DDM failed to capture the shape of RT distributions, although it fit the SAT effect on accuracy
290 and median RTs. The RL-ARD model provides a much better account of the RT distributions,
291 including the differences between SAT conditions. In Figure 6-figure supplement 1 we show
292 that adding non-decision time variability to the RL-DDM mitigates some of the misfit of the
293 RT distributions, although it still consistently under-predicted the 10th percentile in the
294 accuracy condition. Further, this model was still substantially outperformed by the RL-ARD
295 in formal model selection (Δ BPIC = 209), and non-decision time variability was estimated as
296 much greater than what is found in non-learning context, raising the question of its
297 psychological plausibility.

298 Both RL-DDM and RL-ARD models tended to underestimate RTs and choice accuracy in
299 the early trial bins in the accuracy emphasis condition. As in experiment 1, working memory
300 may have contributed to the accurate but slow responses in the first trial bin for the accuracy
301 condition (Collins and Frank, 2018, 2012b; McDougale and Collins, 2020).

302

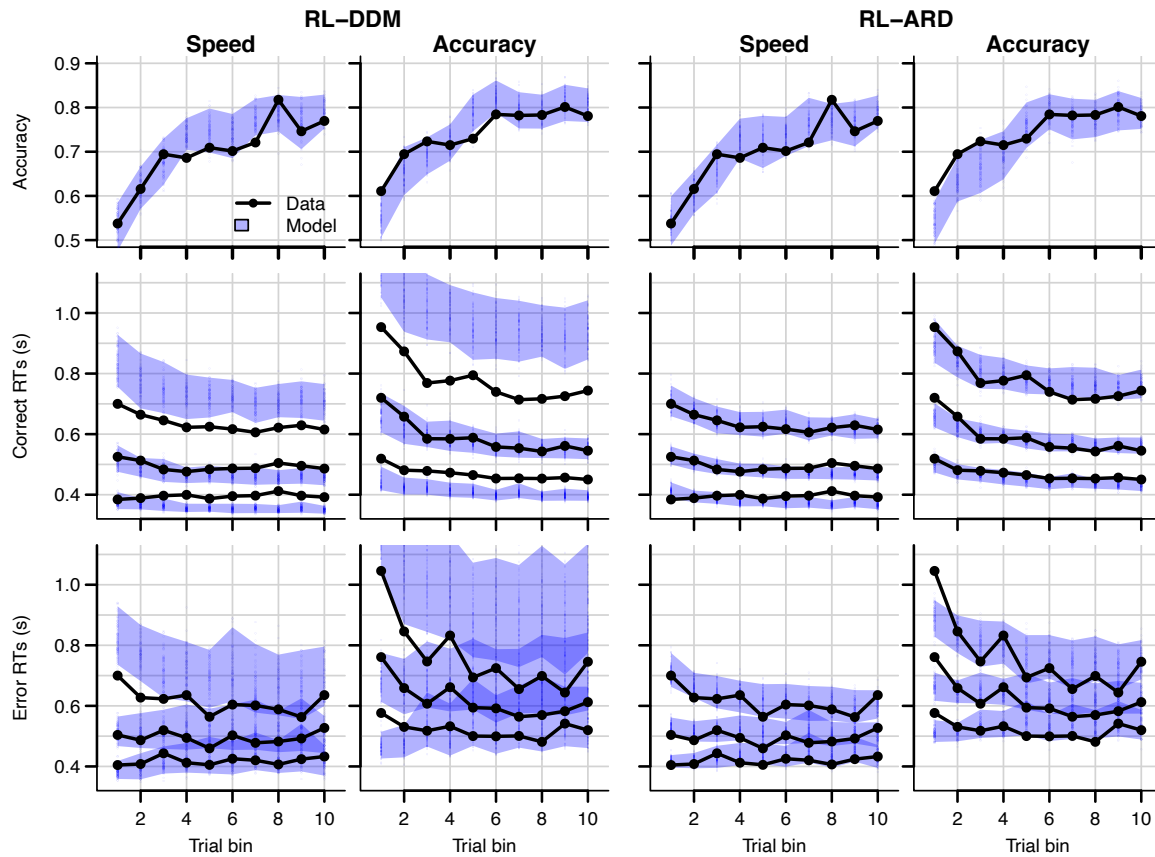


Figure 6. Data (black) and posterior predictive distributions (blue) of the best-fitting RL-DDM (left columns) and the winning RL-ARD model (right columns), separate for the speed and accuracy emphasis conditions. Top row depicts accuracy over trial bins. Middle and bottom row show 10th, 50th, and 90th RT percentiles for the correct (middle row) and error (bottom row) response over trial bins. Shaded areas in the middle and right column correspond to the 95% credible interval of the posterior predictive distribution.

303

304 Reversal learning

305 Finally, we tested whether the RL-ARD can capture changes in accuracy and RTs caused by a
306 perturbation in the learning process due to reversals in reward contingencies. In the reversal
307 learning paradigm (Behrens et al., 2007; Costa et al., 2015; Izquierdo et al., 2017) participants
308 first learn a contingency between choice options and probabilistic rewards (the acquisition
309 phase) that is then suddenly reversed without any warning (the reversal phase). If the link
310 between Q-values and decision mechanisms as proposed by the RL-ARD underlies decisions,
311 the model should be able to account for the behavioral consequences (RT distributions and
312 decisions) of Q-value changes induced by the reversal.

313 Our reversal learning task had the same general structure as experiment 1 (Figure 1), except
314 for the presence of reversals. 47 participants completed four blocks of 128 trials each. Within
315 each block, two pairs of stimuli were randomly interleaved. Between trials 61 and 68
316 (uniformly sampled) in each block, the reward probability switched between stimuli, such that
317 stimuli that were correct during acquisition were incorrect after reversal (and vice versa).
318 Participants were not informed of the reversals prior to the experiment, but many reported
319 noticing them.

320 Data and the posterior predictive distributions of the RL-DDM and the RL-ARD models are
321 shown in Figure 7. Both models captured the change in choice proportions after the reversal
322 reasonably well, although they underestimate the speed of change. In Figure 7-figure
323 supplement 1 we show that the same is true for a standard soft-max model, suggesting that the
324 learning rule is the cause of this problem. Recent evidence indicates that, instead of only
325 estimating expected values of both choice options by error-driven learning, participants may
326 additionally learn the task structure, estimate the probability of a reversal occurring and adjust
327 choice behavior accordingly. Such a model-based learning strategy could increase the speed
328 with which choice behavior changes after a reversal (Costa et al., 2015; Izquierdo et al., 2017;
329 Jang et al., 2015), but as yet a learning rule that implements this strategy has not been
330 developed.

331 The change in RT around the reversal was less marked than the change in choice probability.
332 Once again, the RL-DDM overestimates variability and skew. Both models fit the effects of
333 learning and reversal similarly, but the fastest responses for the RL-DDM decrease much too
334 quickly during initial learning and the reduction in speed for the slowest responses due to the
335 reversal is strongly overestimated. The RL-ARD provides a much better account of the shape
336 of the RT distributions, and furthermore captures the increase in entire RT *distributions*
337 (instead of only the median) after the reversal point. Formal model comparison also very
338 strongly favors the RL-ARD over the RL-DDM ($\Delta BPIC = 4051$). Figure 7-figure supplement
339 2 provides model comparisons to RL-DDMs with between-trial variability parameters, which
340 lead to the same conclusion.

341 A notable aspect of the data is that choice behavior stabilizes approximately 20 trials after
342 the reversal, whereas RTs remain high compared to just prior to the reversal point for up to ~40
343 trials. The RL-ARD explains this behavior through relatively high Q-values for the choice
344 option that was correct during the acquisition (but not reversal) phase (i.e., choice A). Figure
345 8 depicts the evolution of Q-values, Q-value differences and sums, and drift rates in the RL-
346 ARD model. The Q-values for both choice options increase until the reversal (Figure 8A), with
347 a much faster increase for Q_A . At the reversal Q_A decreases and Q_B increases, but as Q_A
348 decreases faster than Q_B increases there is a temporary decrease in Q-value sums (Figure 8C).
349 After approximately 10 trials post-reversal, Q_B is higher than for Q_A , which flips the sign of
350 the Q-value differences (Figure 8B). However, Q_A after the reversal remains higher than the
351 Q_B before the reversal, which causes the (absolute) Q-value differences to be lower after the
352 reversal than before. As a consequence, the drift rates for B after the reversal remain lower than
353 the drift rates for A before the reversal, which increases RT. Clearly, it is important to take
354 account of the sum of inputs to accumulators as well as the difference between them in order
355 to provide an accurate account of the effects of learning.

356

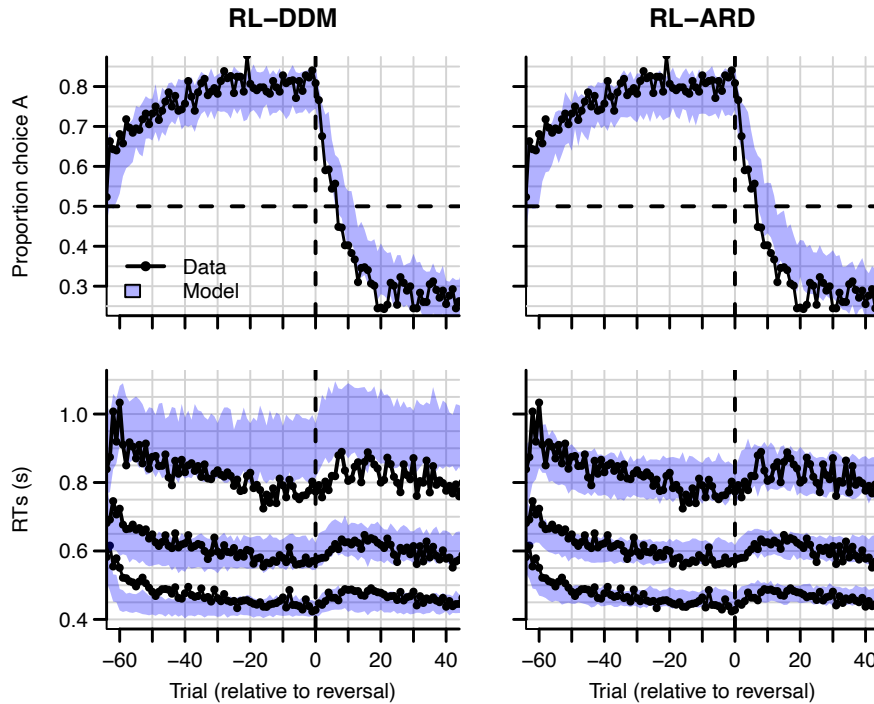


Figure 7. Experiment 3 data (black) and posterior predictive distributions (blue) for the RL-DDM (left) and RL-ARD (right). Top row: choice proportions over trials, with choice option A defined as the high-probability choice before the reversal in reward contingencies. Bottom row: 10th, 50th, and 90th RT percentiles. The data are ordered relative to the trial at which the reversal first occurred (trial 0, with negative trial numbers indicated trials prior to the reversal). Shaded areas correspond to the 95% credible interval of the posterior predictive distributions.

357

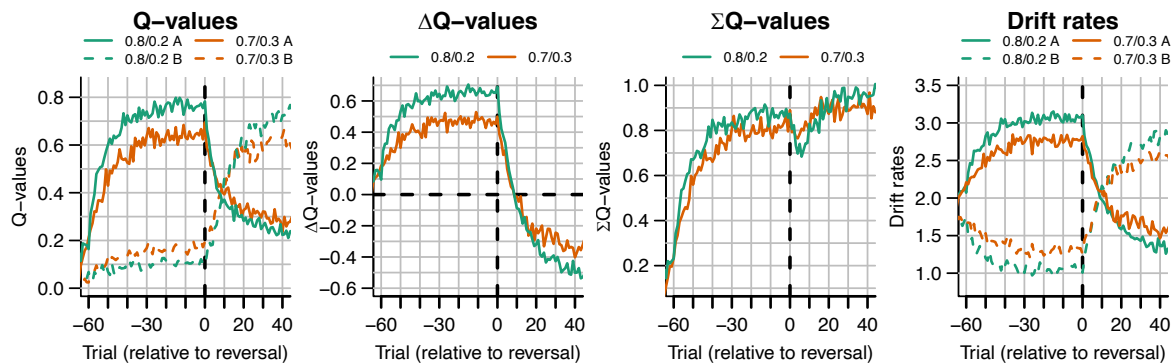


Figure 8. The evolution of Q-values and their effect on drift rates in the RL-ARD in experiment 3, aggregated across participants. Left panel depicts raw Q-values, separate for each difficulty condition (colors). The second and third panel depict the Q-value differences and the Q-value sums over time. The drift rates (right panel) are a weighted sum of the Q-value differences and Q-value sums, plus an intercept. Choice A (solid lines) refers to the option that had the high probability of reward during the acquisition phase, and choice B (dashed lines) to the option that had the high probability of reward after the reversal.

358

359

360 **Discussion**

361 We compared combinations of different evidence-accumulations models with a simple SARSA
362 (Rummery and Niranjan, 1994) reinforcement learning rule (RL-EAMs). The comparison
363 tested the ability of the RL-EAMs to provide a comprehensive account of behavior in learning
364 contexts, not only in terms of the choices made but also the full distribution of the times to
365 make them (RT). We examined a standard instrumental learning paradigm (Frank, 2004) that
366 manipulated the difference in rewards between binary options (i.e., decision difficulty). We
367 also examined two elaborations of that paradigm testing key phenomena from the decision-
368 making and learning literatures, speed-accuracy trade-offs (SAT), and reward reversals,
369 respectively. Our benchmark was the dual threshold Diffusion Decision Model (Ratcliff, 1978)
370 (DDM), which has been used in almost all previous RL-EAM research, but has not been
371 compared to other RL-EAMs, and has not been thoroughly evaluated on its ability to account
372 for RT distributions in learning tasks. Our comparison used several different racing diffusion
373 (RD) models, where decisions depend on the winner of a race between single barrier diffusion
374 processes.

375 The RL-DDM provided a markedly inferior account to the other models, consistently
376 overestimating RT variability and skew. As these aspects of behavior are considered critical in
377 evaluating models in decision-making literature (Forstmann et al., 2016; Ratcliff and McKoon,
378 2008; Voss et al., 2013), our results question whether the RL-DDM provides an adequate
379 model of instrumental learning. Furthermore, the DDM carries with it two important theoretical
380 limitations. First, it can only address binary choice. This is unfortunate given that perhaps the
381 most widely used clinical application of reinforcement learning, the Iowa gambling task
382 (Bechara et al., 1994), requires choices among four options. Second, the input to the DDM
383 combines the evidence for each choice (i.e., “Q” values determined by the learning rule) into a
384 single difference, and so requires extra mechanisms to account for known effects of overall
385 reward magnitude (Fontanesi et al., 2019a). Although there are potential ways that the RL-
386 DDM might be modified to account for magnitude effects, such as increasing between-trial
387 drift rate variability in proportion to the mean rate (Ratcliff et al., 2018), its inability to extend
388 beyond binary choice remains an enduring impediment.

389 The best alternative model that we tested, the RL-ARD (advantage racing diffusion), which
390 is based on the recently proposed advantage accumulation framework (Van Ravenzwaaij et al.,
391 2020), remedied all of these problems. The input to each accumulator is the weighted sum of
392 three components: stimulus independent “urgency”, the difference between evidence for the
393 choice corresponding to the accumulator and the alternative (the advantage), and the sum of
394 the evidence over accumulators. The urgency component had large effect in all fits and played
395 a key role in explaining the effect of speed-accuracy trade-offs. The advantage component,
396 which is similar to the input to the DDM, was strongly supported over a model in which each
397 accumulator only receives evidence favoring its own choice. The sum component provided a
398 simple and theoretically transparent way to deal with reward magnitude effects in instrumental
399 learning. Despite having the weakest effect among the three components, the sum was clearly
400 necessary to provide an accurate fit to our data, even though we did not manipulate reward
401 magnitude. It also played an important role in explaining the effect of reward reversals.

402 It is perhaps surprising that the RL-DDM consistently overestimated RT variability and
403 skewness given that the DDM typically provides much better fits to data from perceptual
404 decision-making tasks without learning. The inclusion of between-trial variability in non-
405 decision times partially mitigated the misfit but required an implausibly high non-decision time
406 variability, and model comparisons still favored the RL-ARD. Previous work on the RL-DDM
407 did not investigate this issue. In many RL-DDM papers, RT distributions are either not
408 visualized at all, or are plotted using (defective) probability density functions on top of a
409 histogram of RT data, making it hard to detect misfit, particularly with respect skew due to the
410 slow tail of the distribution. One exception is Pedersen et al. (2020), whose quantile-based
411 plots show the same pattern that we found here of over-estimated variability and skewness for
412 more difficult choice conditions, despite including between-trial variability in non-decision
413 times. In a non-learning context, it has been shown that the DDM overestimates skewness in
414 high-risk preferential choice data (Dutilh and Rieskamp, 2016). Together these results suggest
415 that decision processes in value-based decision in general, and instrumental learning tasks in
416 particular, may be fundamentally different from a two-sided diffusion process, and instead
417 better captured by a race model such as the RL-ARD.

418 In the current work, we chose to use racing diffusion processes over the more often used
419 LBA models for reasons of parsimony: error-driven learning introduces between-trial
420 variability in accumulation rates, which are explicitly modelled in the RL-EAM framework.
421 As the LBA includes between-trial variability in drift rates as a free parameter, multiple
422 parameters can account for the same variance. Nonetheless, exploratory fits (see Figure 4-
423 figure supplement 3) confirmed our expectation that an RL-ALBA (Advantage ALBA) model
424 fit the data of experiment 1 well, although formal model comparisons preferred the RL-ARD.
425 Future work might consider completely replacing one or more sources of between trial
426 variability in the LBA with structured fluctuations due to learning and adaption mechanisms.

427 The parametrization of the ARD model used in the current paper followed van Ravenzwaaij
428 et al.'s (2020) proposed ALBA model. This parametrization interprets the influence on drift
429 rates in terms of advantages and magnitudes. However, as both the weights on Q-value
430 differences and sums (w_D and w_S) are freely estimated parameters, the equations that define
431 the drift rates can be rearranged as follows:

432

$$\begin{aligned} dx_1 &= [V_0 + w_e Q_1 - w_i Q_2]dt + sW \\ dx_2 &= [V_0 + w_e Q_2 - w_i Q_1]dt + sW \end{aligned} \quad (5)$$

433

434 Where w_e equals the sum of w_D and w_S in the parametrization of Equation (4), and w_i equals
435 the difference between w_D and w_S . This re-parametrization shows that each drift rate is
436 determined by an excitatory influence w_e of the Q-value associated with the accumulator, and
437 an inhibitory influence w_i of the Q-value associated with the other accumulator. Turner (2019)
438 proposed that inhibition plays an important role in learning tasks. Although the locus of
439 inhibition is different in the two models, there are clear parallels that bear further investigation.

440 A limitation of the current work is that we collapsed across blocks in analyzing the data of
441 experiments 2 and 3. However, in more detailed explorations (see Figure 6-figure supplement
442 3) there were indications of second-order changes across blocks. In experiment 2, participants

443 were faster in the first trial bin of the second and third block compared to the first block,
444 suggesting additional practice or adaptation effects at the beginning of the experiment. In
445 experiment 3, participants slowed down, and learned the reversal faster, after the first block.
446 This suggests they learned about the presence of reversals in the first block and applied a
447 different strategy in the later blocks. Although it is known that participants increase their
448 learning rates in volatile environments (Behrens et al., 2007), this by itself does not explain a
449 decrease in response speed. Potentially, if participants understood the task structure after the
450 first block, model-based strategies, such as estimating the probability of a reversal having
451 occurred, also slowed down responses.

452 Although the account of data provided by the RL-ARD model was generally quite accurate,
453 some elements of misfit suggest the need for further model development. RT and accuracy
454 were underestimated in the initial trials of the easiest condition in experiment 1, in the accuracy
455 emphasis condition in experiment 2, and prior to reversals in experiment 3. Furthermore, the
456 RL-ARD model underestimated the speed with which choice probability changed after reversal
457 of stimulus-response mappings. These misfits point to a limited ability to capture the learning-
458 related changes in behavior. This is to some degree unsurprising, since we used a very simple
459 model of error-driven learning. Future work might explore more sophisticated mechanisms,
460 such as multiple learning rates (Daw et al., 2002; Fontanesi et al., 2019a; Gershman, 2015;
461 Pedersen et al., 2017) or different learning rules (Fontanesi et al., 2019b, 2019a). Furthermore,
462 there is clearly a role for working memory in some reinforcement learning tasks (Collins and
463 Frank, 2018, 2012b), likely explaining the accurate but slow responses we observed in the early
464 trial bins for easy conditions.

465 In summary, we believe that the ARD decision mechanism provides a firm basis for further
466 explorations of the mutual benefits that arise from the combination of reinforcement learning
467 and evidence-accumulation models, providing constraint that is based on a more
468 comprehensive account of data than has been possible in the past. As it stands, the RL-ARD's
469 parameter recovery properties are good even with relatively low trial numbers, making it a
470 suitable measurement model for simultaneously studying learning and decision-making
471 processes, and inter-individual differences therein. Further, the advantage framework extends
472 to multiple choice while maintaining analytical tractability and addressing key empirical
473 phenomena in that domain, such as Hick's Law and response-competition effects (Van
474 Ravenzwaaij et al., 2020), enabling future applications to clinical settings, such as in the Iowa
475 gambling task (Bechara et al., 1994).

476

477 **Methods**

478 **Experiment 1**

479 *Participants*

480 61 participants (mean age 21y [SD 2.33], 47 women, 56 right handed) were recruited from the
481 subject pool of the department of Psychology, University of Amsterdam, and participated for
482 course credits. All participants had normal or corrected-to-normal vision and gave written
483 informed consent prior to the experiment onset. The study was approved by the local ethics
484 committee.

485

486 *Task*

487 The task was an instrumental probabilistic learning task (Frank, 2004). On each trial, the
488 subject was presented with two abstract symbols (a ‘stimulus pair’) representing two choice
489 options (see Figure 2A for an example trial). Each choice option had a fixed probability of
490 being rewarded with points when chosen, with one choice option always having a higher
491 probability of being rewarded than the other. The task is to discover, by trial and error, which
492 choice options are most likely to lead to rewards, and thereby to collect as many points as
493 possible.

494 After a short practice block to get familiar with the task, participants completed one block
495 of 208 trials. Four different pairs of abstract symbols were included, each presented 52 times.
496 Stimulus pairs differed in their associated reward probabilities: 0.8/0.2, 0.7/0.3, 0.65/0.35, and
497 0.6/0.4. The size of the reward, if obtained, was always the same: ‘+100’ (or ‘+0’ otherwise).
498 Reward probabilities were chosen such that they differed only in the between-choice difference
499 in reward probability, leading to varying choice difficulties while keeping the mean reward
500 magnitude fixed.

501 Participants were instructed to earn as many points as possible, and to always respond before
502 the deadline of 2 seconds. Feedback consisted of two parts: an ‘outcome’ and a ‘reward’. The
503 outcome corresponded to the probabilistic outcome of the choice, whereas the reward
504 corresponded to the actual number of earned points. When participants responded before the
505 deadline, the reward was equal to the outcome. If they were too late, the outcome was shown
506 to allow participants to learn from their choice, but the reward they received was set to 0 to
507 encourage responding in time. Participants received a bonus depending on the number of points
508 earned (maximum +0.5 course credits, mean received +0.24). The task was coded in PsychoPy
509 (Peirce et al., 2019). After this block, participants performed two more blocks of the same task
510 with different manipulations, which are not of current interest.

511 512 *Exclusion*

513 Six participants were excluded from analysis: One reported, after the experiment, not to have
514 understood the task, one reported a technical issue, and four did not reach an above-chance
515 accuracy level as determined by a binomial test (accuracy cut-off 0.55, corresponding to
516 $p < 0.05$). The final sample thus consisted of 55 subjects (14 men, mean age 21 years old [SD
517 2.39], 51 right-handed).

518 519 *Cognitive modelling*

520 The main analysis consists of fitting four RL-EAMs to the data and comparing the quality of
521 the fits penalized by model complexity. We compared four different decision models: the DDM
522 (Ratcliff, 1978), a racing diffusion (Boucher et al., 2007; Logan et al., 2014; Purcell et al.,
523 2010; Turner, 2019) model, and two Advantage Racing Diffusion (ARD; Van Ravenzwaaij et
524 al., 2020) models (see Figure 1 for an overview). Whereas the former is a two-sided diffusion
525 process, the latter three models employ a race architecture.

526 For all models we used the State-Action-Reward-State-Action (SARSA; Rummery and
527 Niranjana, 1994) update rule as a learning model:

$$528 \quad Q_{i,t+1} = Q_{i,t} + \alpha(r_t - Q_{i,t}) \quad (4)$$

529

530 where $Q_{i,t}$ is the value representation of choice option i on trial t , α the learning rate, and r_t
531 the reward on trial t . The difference between the actual reward and the value representation of
532 the chosen stimulus, $r_t - Q_{i,t}$, is known as the reward prediction error. The learning rate
533 controls the speed at which Q-values change in response to the reward prediction error, with
534 larger learning rates leading to stronger fluctuations. In this model, only the Q-value of the
535 chosen option is updated.

536

537 *RL-EAM 1: RL-DDM*

538 In the first RL-EAM, we use the DDM (Ratcliff, 1978) as a choice model (Figure 1, left
539 column). The DDM assumes that evidence accumulation is governed by:

540

$$dx = vdt + sW$$

541

542 v is the mean speed of evidence accumulation (the *drift rate*), and s is the standard deviation
543 of the within-trial accumulation white noise (W). The RL-DDM assumes that the drift rate
544 depends linearly on the difference of value representations:

545

$$v_t = w(Q_{t,1} - Q_{t,2})$$

546

547 w is a weighting variable, and $Q_{t,1}$ and $Q_{t,2}$ are the Q-values for both choice options per trial,
548 which change each trial according to Equation 4. Hence,

549

$$dx = w(Q_1 - Q_2) dt + sW \quad (1)$$

550

551 The starting point of evidence accumulation, z , lies between decision boundaries a and $-a$.
552 Here, as in earlier RL-DDM work (Fontanesi et al., 2019a, 2019b; Pedersen et al., 2017), we
553 assume an unbiased start of the decision process (i.e., $z = 0$). Evidence accumulation finishes
554 when threshold a or $-a$ is reached, and the decision for the choice corresponding to Q_1 or Q_2 ,
555 respectively, is made. The response time is the time required for the evidence-accumulation
556 process to reach the bound, plus an intercept called the non-decision time (t_0). The non-
557 decision time is the sum of the time required for perceptual encoding and the time required for
558 the execution of the motor response. Parameter s was fixed to 1 to satisfy scaling constraints
559 (Donkin et al., 2009; van Maanen and Miletic, 2020). In total, this specification of the RL-
560 DDM has 4 free parameters (α, w, a, t_0).

561 Furthermore, we fit four additional RL-DDMs (RL-DDM A1-A4) with between-trial
562 variabilities in start point, drift rate, and non-decision time, as well as a non-linear link function
563 between Q-values and drift rates (Fontanesi et al., 2019a). RL-DDM A1 uses the non-linear
564 function $v_t = \frac{2v_{max}}{1 + \exp(w(Q_{t,1} - Q_{t,2}))} - v_{max}$ to link Q-values to drift rates (Fontanesi et al., 2019a).

565 For the new v_{max} parameter, $\mathcal{N}(2,5)$ (truncated 0) and $\Gamma(1,1)$ were used as priors for the
566 hypermean and hyperSD, respectively. RL-DDM A2 includes between-trial variabilities in
567 both drift rate s_v and start point s_z , with $\mathcal{N}(0.1,0.1)$ and $\mathcal{N}(0.1,0.1)$ as priors for hypermeans
568 (respectively, both truncated at 0) and $\Gamma(1,1)$ for the hyperSD. Drift rate variability was

569 estimated as a proportion of the current drift rate, such that $s_{v,t} = v_t * s_v$ (which allows for
570 higher variability terms for higher Q-value differences, but retains the ratio v/s_v). RL-DDM
571 A3 included s_v , s_z , and also between-trial variability in non-decision time s_{t0} , for which
572 $\mathcal{N}(0.1,0.1)$ (truncated at 0) and $\Gamma(1,1)$ were used as priors for the hypermean and hyperSD,
573 respectively. RL-DDM A4 used all three between-trial variabilities as well as the non-linear
574 link function. The quality of fits of these additional models can be found in Figure 3-figure
575 supplement 1. Foreshadowing the results, the RL-DDM A3 improved the quality of fit
576 compared to the RL-DDM, but required an implausibly high non-decision time variability: The
577 across-subject mean of the median posterior estimates of the $t0$ and s_{t0} parameters indicate a
578 non-decision time distribution of [0.27 s, 0.64 s]. The range of 0.37 s is very high in light of
579 the literature (Tran et al., n.d.), raising the question of its psychological plausibility. For this
580 reason, as well as since the RL-DDM is used most often without s_{t0} , we focus on the RL-DDM
581 (without between-trial variabilities) in the main text.

582

583 *RL-EAM 2: RL-RD*

584 The RL-RD (Figure 1, middle panel) assumes that two evidence accumulators independently
585 accrue evidence for one choice option each, both racing towards a common threshold a
586 (assuming no response bias). The first accumulator to hit the bound wins, and the
587 corresponding decision is made. For each choice option i , the dynamics of accumulation are
588 governed by:

589

$$dx_i = [V_0 + wQ_i]dt + sW \quad (2)$$

590

591 V_0 is a parameter specifying the drift rate in the absence of any evidence, w a weighting
592 parameter, and s the standard deviation of within-trial noise. As such, the mean speed of
593 accumulation (the drift rate v_i) is the sum of two independent factors: an evidence-independent
594 baseline speed V_0 , and an evidence-dependent weighted Q-value, wQ_i . Since V_0 is assumed to
595 be identical across accumulators, and governs the speed of accumulation unrelated to the
596 amount of evidence, we interpret this parameter as an additive urgency signal (Miletić and Van
597 Maanen, 2019), with conceptually similar behavioral effects as collapsing bounds (Hawkins et
598 al., 2015). Similar to the DDM, a non-decision time parameter accounts for the time for
599 perceptual encoding and the motor response time. Parameter s was fixed to 1 to satisfy scaling
600 constraints (Donkin et al., 2009; van Maanen and Miletić, 2020). In total, the RL-RD has 5 free
601 parameters ($\alpha, w, a, v0, t0$).

602 Each accumulator's first passage times are Wald (also known as inverted Gaussian)
603 distributed (Anders et al., 2016). In an independent race model, each accumulator's first
604 passage time distribution is normalized to the probability of the response with which it is
605 associated (Brown and Heathcote, 2008; Turner, 2019).

606

607 *RL-EAM 3 & 4: RL-ARD*

608 Thirdly, we fit two racing diffusion models based on an advantage race architecture (Van
609 Ravenzwaaij et al., 2020). An advantage race model using an LBA has been shown to provide
610 a natural account for multi-alternative choice phenomena such as Hick's law, as well as

611 stimulus magnitude effects in perceptual decision-making. Like in the RL-RD, accumulators
612 race towards a common bound, but the speed of evidence accumulation v_i depends on multiple
613 factors: first, as in the RL-RD, the evidence-independent speed of accumulation V_0 ; second,
614 the *advantage* of the evidence for one choice option over the other (c.f. the DDM, where the
615 difference between evidence for both choice options is accumulated); and third, the *sum* of the
616 total available evidence. Combined, for two accumulators in the RL-EAM framework, this
617 leads to:

$$\begin{aligned} dx_1 &= [V_0 + w_d(Q_1 - Q_2) + w_s(Q_1 + Q_2)]dt + sW \\ dx_2 &= [V_0 + w_d(Q_2 - Q_1) + w_s(Q_1 + Q_2)]dt + sW \end{aligned} \quad (3)$$

619
620 In the original work proposing the advantage accumulation framework (Van Ravenzwaaij et
621 al., 2020), it was shown that the w_d parameter had a much stronger influence on evidence-
622 accumulation rates than the w_s parameter. Therefore, we first fixed the w_s parameter to 0, to
623 test whether the accumulation of *differences* is sufficient to capture all trends in the data. We
624 term this model the RL-lARD (l = limited), which we compare to the RL-ARD in which we fit
625 w_s as a free parameter.

626 As previously, parameter s was fixed to 1 to satisfy scaling constraints (Donkin et al., 2009;
627 van Maanen and Miletic, 2020). The RL-ARD also has a threshold, non-decision time, and
628 learning rate parameter, totaling five (α, w_d, a, V_0, t_0) and 6 free parameters
629 ($\alpha, w_d, w_s, a, V_0, t_0$) for the RL-lARD and RL-ARD, respectively. A parameter recovery study
630 was (Heathcote et al., 2015; Miletic et al., 2017; Moran, 2016; Spektor and Kellen, 2018)
631 performed to confirm that data-generating parameters can be recovered using the experimental
632 paradigm at hand. The results are shown in Figure 3-figure supplement 2.

633
634 *Bayesian hierarchical parameter estimation, posterior predictive distributions, model*
635 *comparisons*

636 We estimated group-level and subject-level posterior distributions of each model's parameter
637 using a combination of differential evolution (DE) and Markov-chain Monte Carlo sampling
638 (MCMC) with Metropolis-Hastings (Ter Braak, 2006; Turner et al., 2013). Sampling settings
639 were default as implemented in the Dynamic Models of Choice *R* software (Heathcote et al.,
640 2019): The number of chains, D , was three times the number of free parameters. Cross-over
641 probability was set to $2.38/\sqrt{D}$ at the subject-level and $U[0, 1]$ at the group level. Migration
642 probability was set to 0.05 during burn-in only. Convergence was assessed using visual
643 inspection of the chain traces and Gelman-Rubin diagnostic (Brooks and Gelman, 1998;
644 Gelman and Rubin, 1992) (individual and multivariate potential scale factors < 1.03 in all
645 cases).

646 Hierarchical models were fit assuming independent normal population ("hyper")
647 distributions for each parameter. For all models, we estimated the learning rate on a probit scale
648 (mapping $[0, 1]$ onto the real domain), with a normal prior $\alpha \sim \Phi(\mathcal{N}(-1.6, 5))$ (Spektor and
649 Kellen, 2018). Prior distributions for all estimated hyper-mean decision-related parameters
650 were vague. RL-EAMs, the threshold parameter $a \sim \mathcal{N}(3, 5)$ truncated at 0, and
651 $t_0 \sim \mathcal{N}(0.3, 0.5)$ truncated at 0.025 s and 1 s (all estimation was carried out on the seconds

652 scale). For the RL-DDM, $w \sim \mathcal{N}(2, 5)$. For the RL-RD, $w \sim \mathcal{N}(9, 5)$, and for the RL-ARD
653 models, $w_D \sim \mathcal{N}(9, 5)$ and $w_S \sim \mathcal{N}(0, 3)$. For the hyper-SD, a $\Gamma(1,1)$ distribution was used
654 as prior. Plots of superimposed prior and posterior hyper-distributions confirmed that these
655 prior setting were not influential.

656 In initial explorations, we also freely estimated the Q-values at trial 0. However, in the RL-
657 EAMs, the posterior distributions for these Q-values consistently converged on 0, which was
658 therefore subsequently used as a fixed value for all results reported here. For the soft-max fits,
659 they were set to 0.5 as often used in reinforcement learning models of two-choice tasks (Apps
660 et al., 2015; Collins and Frank, 2018, 2012b; Fontanesi et al., 2019a; McDougale and Collins,
661 2020; Pedersen and Frank, 2020). Including the initial Q-values as a free parameter in the soft-
662 max models of experiment 2 led to the same conclusions.

663 To visualize the quality of model fit, we took 100 random samples from the estimated
664 parameter posteriors and simulated the experimental design with these parameters. For each
665 behavioral measure (e.g., RT quantiles, accuracy), credible intervals were estimated by taking
666 the range between the 2.5% and 97.5% quantiles of the averages over participants.

667 To quantitatively compare the fit of different models, penalized by their complexity, we
668 used the Bayesian predictive information criterion (BPIC; Ando, 2007). The BPIC is an
669 analogue of the Bayesian information criterion (BIC), but (unlike the BIC) suitable for models
670 estimated using Bayesian methods. Compared to the deviance information criterion (DIC;
671 Spiegelhalter et al., 2002), the BPIC penalizes model complexity more strongly to prevent
672 over-fitting (c.f. AIC vs. BIC). Lower BPIC values indicate better trade-offs between fit quality
673 and model complexity.

674

675 **Experiment 2**

676 *Participants*

677 23 participants (mean age 19 years old [SD 1.06 years], 7 men, 23 right-handed) were recruited
678 from the subject pool of the Department of Psychology of the University of Amsterdam and
679 participated for course credits. Participants did not participate in experiment 1 or 3. All
680 participants had normal or corrected-to-normal vision and gave written informed consent prior
681 to the experiment onset. The study was approved by the local ethics committee.

682

683 *Task*

684 Participants performed the same task as in experiment 1, with the addition of an SAT
685 manipulation (Figure 2C). The SAT manipulation included both an instructional cue and a
686 response deadline. Prior to each trial, a cue instructed participants to emphasize either decision
687 speed ('SPD') or decision accuracy ('ACC') in the upcoming trial, and in speed trials,
688 participants did not earn points if they were too late (> 700 ms). As in experiment 1, after each
689 choice participants received feedback consisting of two components: an outcome and a reward.
690 The outcome refers to the outcome of the probabilistic gamble, whereas the reward refers to
691 the number of points participants actually received. If participants responded in time, the
692 reward was equal to the outcome. In speed trials, participants did not earn points if they
693 responded later than 700 ms after stimulus onset, even if the outcome was +100. On trials
694 where participants responded too late, they were additionally informed of the reward that was

695 associated with their choice, had they been in time. This way, even when participants are too
696 late, they still receive the feedback that can be used to learn from their choices.

697 The deadline manipulation was added because we hypothesized that instructional cues alone
698 would not be sufficient to persuade participants to change their behavior in the instrumental
699 learning task, since that task specifically requires them to accumulate points. If the received
700 number of points was independent of response times, the optimal strategy to collect most points
701 would be to ignore the cue and focus on accuracy only.

702 Participants performed 324 trials divided over 3 blocks. Within each block, three pairs of
703 stimuli were shown, with associated reward probabilities of 0.8/0.2, 0.7/0.3, and 0.6/0.4. Speed
704 and accuracy trials were randomly interleaved. Figure 2C depicts the sequence of events in
705 each trial. As this experiment also served as a pilot for an fMRI experiment, we added fixation
706 crosses between each phase of the trial, with jittered durations. A pre-stimulus fixation cross
707 lasted 0.5, 1, 1.5, or 2 s; fixation crosses between cue and stimulus, between stimulus and
708 highlight, and between highlight and feedback lasted 0, 0.5, 1, or 1.5 s; and an inter-trial
709 interval fixation cross lasted 0.5, 1, 1.5, 2, 2.5 seconds. Each trial took 7.5 seconds. The
710 experiment took approximately 45 minutes.

711

712 *Exclusion*

713 Four participants did not reach above-chance performance as indicated by a binomial test (cut-
714 off 0.55, $p < 0.05$), and were excluded from further analyses. The final sample thus consisted
715 on 19 participants (mean age 19 years old [SD 1.16 years], 6 men, 19 right-handed). For one
716 additional participant, a technical error occurred after the first block. This participant was
717 included in the analyses, since the Bayesian estimation framework naturally down-weighs the
718 influence of participants with fewer trials.

719

720 *Manipulation check & across-block differences in behavior*

721 We expected an interaction between SAT conditions and learning. In the early trials,
722 participants have not yet learned the reward contingencies, causing a low evidence
723 accumulation rate compared to later trials. With low rates it takes longer to reach the decision
724 threshold, and small changes in the threshold settings or drift rates (by means of an additive
725 urgency signal) can cause large behavioral effects. Therefore, we expected the behavioral
726 effects of the SAT manipulation to become smaller over the course of learning.

727 To formally test for the behavioral effects of the SAT manipulation in experiment 2, we fit
728 two mixed effects models (Gelman and Hill, 2007): A linear model with RT as dependent
729 variable, and a logistic model with accuracy as the dependent variable. As fixed effects, trial
730 bin and SAT condition were included. Trial bins were obtained by splitting all trials in ten bins
731 (approximately 20 trials each) per participants. As random effects, only participant was
732 included. For the logistic model using accuracy as a dependent variable, we log-transformed
733 trial bin numbers (Evans et al., 2018; Heathcote et al., 2000), to account for the non-linear
734 relation between accuracy and trial bin (Figure 6, top row). Mixed effects analyses were done
735 using *lme4* (Bates et al., 2015). For all mixed effects models, we report parameter estimates of
736 the fixed effects, their standard error and confidence interval, as well as a p -value obtained
737 from a t -distribution with the denominator degrees of freedom approximated using

738 Satterthwaite's method (Satterthwaite, 1941), as implemented in the *lmerTest* package
739 (Kuznetsova et al., 2017) for the *R* programming language (R Core Team, 2017).

740 Next, we tested for the across block stability of behavior using two mixed effects models.
741 One linear mixed effects model was used to predict RT with block number, trial bin, and their
742 interaction, with a random intercept for participant. A second, logistic mixed effects model was
743 used to test the effect of block, trial bin (log-transformed, as above), and their interaction on
744 choice accuracy. In both models, trial bin is expected to influence the outcome variables, but
745 the assumption of across block stability in behavior is violated if there are main effects of block
746 number and/or interaction effects between block number and trial bin. Mean RT and accuracy
747 by block is shown together with the formal test results in Figure 6-figure supplement 3.

748

749 *Cognitive modelling*

750 First, we tested whether a standard soft-max model is able to capture the difference in choice
751 behavior. Soft-max is given by:

752

$$P_{i,t} = \frac{\exp \beta Q_{i,t}}{\sum_j \exp \beta Q_{j,t}} \quad (5)$$

753

754 where $P_{i,t}$ is the probability of choosing option i on trial t , J is the total number of choice
755 options, and β is a free parameter often called the inverse temperature. The inverse temperature
756 is often interpreted in terms of the exploration/exploitation trade-off (Daw et al., 2006), with
757 higher values indicating more exploitation. In two-choice settings, Equation 5 can be re-written
758 as:

759

$$P_{2,t} = \frac{1}{1 + \exp \beta (Q_{1,t} - Q_{2,t})} \quad (6)$$

760

761 which highlights that the choice probability is driven by the *difference* in Q-values, weighted
762 by the inverse temperature parameter. We hierarchically fit two soft-max models using the
763 same parameter estimation methods as in experiment 1. One model assumed a single β
764 parameter, the other model assumed a β parameter per SAT condition. Priors for the hypermean
765 were set to $\beta \sim N(1,5)$ truncated at 0, and for the hyperSD $\Gamma(1,1)$.

766 Next, we fit three RL-DDMs and seven RL-ARDs. The three RL-DDM models varied either
767 threshold, the Q-value weighting on the drift rates parameter (Sewell and Stallman, 2020), or
768 both. The seven RL-ARD allowed all unique combinations of the threshold, urgency, and drift
769 rate parameters free to vary between the speed and accuracy conditions.

770 For the accuracy condition, we used the same priors as in experiment 1. In the speed
771 condition, the parameters that were free to vary were estimated as proportional differences
772 from the accuracy conditions; specifically: $a_{spd} = (1 + m_{a,spd}) * a_{acc}$, $V_{0,spd} = (1 +$
773 $m_{V_0,spd}) * V_{0,acc}$, and $v_{i,spd} = (1 + m_{v,spd}) * v_{i,acc}$. The prior used was $\mathcal{N}(0,5)$ for the
774 hypermean and $\Gamma(1,1)$ for the hyperSD of all parameters m , truncated at -1.

775 As in experiment 1, we performed a parameter recovery study to confirm that the data-
776 generating parameters can be recovered, using the winning model and a simulation of the
777 paradigm of experiment 2. The results are shown in Figure 6-figure supplement 2.

778 Additionally, we performed a second model comparison using three variants of RL-DDM
779 A3 (i.e., including between-trial variabilities in start point, drift rate, and non-decision time),
780 which varied the threshold, the Q-value weighting on the drift rate parameters, and both. The
781 results are shown in Figure 6-figure supplement 1, and lead to the same conclusions as the RL-
782 DDM.

783

784 **Experiment 3**

785 *Participants*

786 47 participants (mean age 21 years old [SD 2.81 years], 16 men, 40 right-handed) were
787 recruited from the subject pool of the Department of Psychology of the University of
788 Amsterdam and participated for course credits. Participants did not participate in experiment 1
789 or 2. All participants had normal or corrected-to-normal vision and gave written informed
790 consent prior to the experiment onset. The study was approved by the local ethics committee.

791

792 *Task*

793 The reversal learning task had the same general task structure as experiment 1. Participants
794 completed four blocks of 128 trials each, totaling 512 trials. Within each block, two pairs of
795 stimuli were randomly interleaved, with associated reward probabilities of 0.8/0.2 and 0.7/0.3.
796 Between trials 61 and 68 (uniformly sampled) of each block, the reward probability switched
797 between stimuli, such that the stimulus with a pre-reversal reward probability of 0.8/0.7 had a
798 post-reversal reward probability of 0.2/0.3 (and vice versa). Participants were not informed of
799 the reversals prior to the experiment, but many reported noticing them.

800 In addition to the reversal learning task, the experimental session also contained a working
801 memory task that is not of current interest. 30 participants performed the reversal learning task
802 before the working memory task, and 17 participants afterwards. The entire experiment took
803 approximately one hour.

804

805 *Cognitive modelling*

806 The RL-DDM and RL-ARD were fit to the data using the same methods as in experiment 1.
807 Again, we performed a parameter recovery study, of which the results are shown in Figure 7-
808 figure supplement 3. Similarly, we also fit RL-DDM A3 to the data. In an initial fit, the MCMC
809 chains for 11 (out of 47) participants got stuck in values for s_z of 1 (i.e., s_z covered the entire
810 range between both thresholds), which are implausibly high and moreover led to convergence
811 problems. We re-fit this model with the prior on $s_z \sim \mathcal{N}(0.1, 0.1)$ truncated to 0 and 0.5 (i.e.,
812 setting the maximum range of between-trial start point variability to be half the range between
813 the lower and upper threshold), which did converge. The posterior predictives are shown in
814 Figure 7-figure supplement 2). This model led to the same overall conclusions as the standard
815 RL-DDM.

816

817

818 Table 1. Posterior parameter estimates (across-subject mean and SD of the median of the
 819 posterior distributions) for all models and experiments. For models including s_{t0} , the non-
 820 decision time is assumed to be uniformly distributed with bounds $[t0, t0 + s_{t0}]$.
 821

Experiment 1									
RL-DDM	α	a	$t0$	w					BPIC
	0.14 (0.11)	1.48 (0.19)	0.30 (0.06)	3.21 (1.11)					7673
RL-RD	α	a	$t0$	V_0	w				
	0.12 (0.08)	2.16 (0.27)	0.10 (0.04)	1.92 (0.42)	3.09 (1.32)				5613
RL-IARD	α	a	$t0$	V_0	w_d				
	0.13 (0.12)	2.05 (0.24)	0.12 (0.05)	2.48 (0.43)	2.36 (0.95)				4849
RL-ARD	α	a	$t0$	V_0	w_d	w_s			
	0.13 (0.11)	2.14 (0.26)	0.11 (0.04)	2.46 (0.59)	2.25 (0.78)	0.36 (0.79)			4577
RL-DDM A1	α	a	$t0$	w	v_{max}				
	0.14 (0.12)	1.49 (0.20)	0.30 (0.06)	3.01 (0.66)	2.81 (0.72)				7717
RL-DDM A2	α	a	$t0$	w	s_z	s_v			
	0.14 (0.11)	1.48 (0.19)	0.30 (0.06)	3.21 (1.12)	$1.79e^{-3}$ ($0.4e^{-3}$)	$1.8e^{-3}$ ($0.4e^{-3}$)			7637
RL-DDM A3	α	a	$t0$	w	s_z	s_v	s_{t0}		
	0.13 (0.12)	1.13 (0.19)	0.27 (0.06)	5.31 (2.04)	0.00 (0.00)	0.31 (0.13)	0.37 (0.13)		4844
RL-DDM A4	α	a	$t0$	w	v_{max}	s_v	s_z	s_{t0}	
	0.13 (0.12)	1.15 (0.17)	0.27 (0.06)	2.02 (0)	5.16 (1.18)	0.55 (0.24)	$1.57e^{-3}$ (0)	0.36 (0.13)	4884
RL-ALBA	α	a	$t0$	V_0	w_d	w_s	A		
	0.13 (0.11)	3.53 (0.53)	0.03 (0.00)	3.03 (0.57)	2.03 (0.59)	0.33 (0.78)	1.73 (0.43)		4836
Experiment 2									
Soft-max 1	α	β							
	0.08 (0.05)	7.29 (2.59)							6278
Soft-max 2	α	β_{spd}	β_{acc}						
	0.08 (0.05)	6.51 (1.49)	8.67 (2.02)						6289
RL-DDM 1	α	$a_{spd} /$ a_{acc}	$t0$	w					
	0.13 (0.06)	1.11 (0.18) / 1.42 (0.23)	0.26 (0.06)	3.28 (0.66)					979
RL-DDM 2	α	a	$t0$	w_{spd}/w_{acc}					
	0.13 (0.05)	3.01 (0.63)	0.26 (0.06)	3.46 (0.79) / 3.01 (0.63)					1518
RL-DDM 3	α	$a_{spd} /$ a_{acc}	$t0$	w_{spd}/w_{acc}					
	0.13 (0.06)	1.10 (0.18) / 1.44 (0.23)	0.26 (0.06)	3.11 (0.68) / 3.48 (0.72)					999
RL-ARD 1	α	a_{spd} / a_{acc}	$t0$	V_0	w_d	w_s			

	0.12 (0.05)	1.45 (0.35) / 1.82 (0.35)	0.15 (0.07)	2.59 (0.50)	2.24 (0.53)	0.47 (0.34)		-1044
RL-ARD 2	α	a	t_0	V_0	w_d	w_s	$m_{v,spd}$	
	0.12 (0.05)	1.83 (0.36)	0.12 (0.07)	2.52 (0.53)	1.83 (0.56)	0.32 (0.26)	1.31 (0.20)	-827
RL-ARD 3	α	a	t_0	$V_{0,spd}/$ $V_{0,acc}$	w_d	w_s		
	0.12 (0.05)	1.83 (0.35)	0.12 (0.07)	3.37 (0.84) / 3.37 (0.54)	2.11 (0.52)	0.39 (0.30)		-934
RL-ARD 4	α	a_{spd} / a_{acc}	t_0	V_0	w_d	w_s	$m_{v,spd}$	
	0.12 (0.05)	1.04 (0.14) / 1.82 (0.35)	0.15 (0.07)	2.59 (0.52)	2.21 (0.51)	0.44 (0.38)	1.04 (0.14)	-1055
RL-ARD 5	α	$a_{spd}/$ a_{acc}	t_0	$V_{0,spd}/$ $V_{0,acc}$	w_d	w_s		
	0.12 (0.05)	1.59 (0.40) / 1.83 (0.32)	0.14 (0.06)	2.92 (0.65) / 2.52 (0.50)	2.21 (0.50)	0.43 (0.33)		-1071
RL-ARD 6	α	a	t_0	$V_{0,spd}/$ $V_{0,acc}$	w_d	w_s	$m_{v,spd}$	
	0.12 (0.05)	1.86 (0.35)	0.12 (0.07)	4.13 (0.98) / 2.40 (0.54)	2.28 (0.53)	0.44 (0.33)	0.84 (0.03)	-897
RL-ARD 7	α	$a_{spd}/$ a_{acc}	t_0	$V_{0,spd}/$ $V_{0,acc}$	w_d	w_s	$m_{v,spd}$	
	0.12 (0.05)	1.61 (0.40) / 1.87 (0.32)	0.14 (0.06)	3.66 (0.74) / 2.52 (0.50)	2.41 (0.53)	0.48 (0.38)	0.82 (0.08)	-1060
RL-DDM A3 1	α	$a_{spd} /$ a_{acc}	t_0	w	s_z	s_v	s_{t_0}	
	0.12 (0.05)	0.81 (0.16) / 1.14 (0.17)	0.23 (0.06)	4.46 (0.79)	0.10 (0.01)	0.18 (0.05)	0.26 (0.09)	-862
RL-DDM A3 2	α	a	t_0	w_{spd}/w_{acc}	s_z	s_v	s_{t_0}	
	0.12 (0.05)	1.03 (0.14)	0.24 (0.06)	18.4 (23.34) / 4.44 (0.84)	0.26 (0.07)	0.61 (0.50)	0.28 (0.10)	-325
RL-DDM A3 3	α	$a_{spd} /$ a_{acc}	t_0	w_{spd}/w_{acc}	s_z	s_v	s_{t_0}	
	0.12 (0.05)	0.81 (0.16) / 1.14 (0.17)	0.23 (0.06)	4.45 (0.83) / 4.45 (0.83)	0.07 (0.00)	0.17 (0.04)	0.26 (0.09)	-849
Experiment 3								
Soft-max	α	β						
	0.36 (0.17)	3.5 (1.83)						24568
RL-DDM	α	a	t_0					

	0.38 (0.14)	1.37 (0.24)	0.24 (0.07)					15599
RL-ARD	α	a	t_0	V_0	w_d	w_s		
	0.35 (0.15)	1.48 (0.34)	0.13 (0.08)	1.86 (0.51)	1.52 (0.63)	0.23 (0.25)		11548
RL-DDM A3	α	a	t_0	w	s_z	s_v	$s_{t,0}$	
	0.38 (0.14)	1.15 (0.22)	0.22 (0.07)	2.72 (1.16)	0.21 (0.09)	0.28 (0.15)	0.27 (0.17)	11659

822

823

824

825

826 Data availability statement

827 All data are available on OSF (<https://osf.io/ygrve/>).

828

829 Code availability statement

830 All analysis code is available on OSF (<https://osf.io/ygrve/>).

831

832 Acknowledgements

833 We thank Barbara Mathiopoulou and Chris Riddell for their help collecting the data. This work
834 was supported by an NWO-VICI grant (BUF), an ABC VIP grant and ARC DP150100272 and
835 DP160101891 grants (AH).

836

837

838

839 References

- 840 Anders R, Alario F, Van Maanen L. 2016. The Shifted Wald Distribution for Response Time Data Analysis.
841 *Psychol Methods* **21**:309–327.
- 842 Ando T. 2007. Bayesian predictive information criterion for the evaluation of hierarchical Bayesian and
843 empirical Bayes models. *Biometrika* **94**:443–458. doi:10.1093/biomet/asm017
- 844 Apps MAJ, Lesage E, Ramnani N. 2015. Vicarious reinforcement learning signáis when instructing others. *J*
845 *Neurosci* **35**:2904–2913. doi:10.1523/JNEUROSCI.3669-14.2015
- 846 Arnold NR, Bröder A, Bayen UJ. 2015. Empirical validation of the diffusion model for recognition memory and
847 a comparison of parameter-estimation methods. *Psychol Res* **79**:882–898. doi:10.1007/s00426-014-0608-y
- 848 Barto AG, Sutton RS, Brouwer PS. 1981. Associative search network: A reinforcement learning associative
849 memory. *Biol Cybern* **40**:201–211. doi:10.1007/BF00453370
- 850 Bates D, Mächler M, Bolker B, Walker S. 2015. Fitting Linear Mixed-Effects Models Using lme4. *J Stat Softw*
851 **67**. doi:10.18637/jss.v067.i01
- 852 Bechara a, Damasio a R, Damasio H, Anderson SW. 1994. Insensitivity to future consequences following
853 damage to human prefrontal cortex. *Cognition* **50**:7–15. doi:10.1016/0010-0277(94)90018-3
- 854 Behrens TEJ, Woolrich MW, Walton ME, Rushworth MFS. 2007. Learning the value of information in an
855 uncertain world. *Nat Neurosci* **10**:1214–1221. doi:10.1038/nn1954
- 856 Boag RJ, Strickland L, Heathcote A, Neal A, Loft S. 2019a. Cognitive Control and Capacity for Prospective
857 Memory in Complex Dynamic Environments. *J Exp Psychol Gen* **148**:2181–2206.
858 doi:10.1037/xge0000599
- 859 Boag RJ, Strickland L, Loft S, Heathcote A. 2019b. Strategic attention and decision control support prospective
860 memory in a complex dual-task environment. *Cognition* **191**:103974. doi:10.1016/j.cognition.2019.05.011
- 861 Bogacz R, Larsen T. 2011. Integration of reinforcement learning and optimal decision-making theories of the
862 basal ganglia. *Neural Comput* **23**:817–851. doi:10.1162/NECO_a_00103
- 863 Bogacz R, McClure SM, Li J, Cohen JD, Montague PR. 2007. Short-term memory traces for action bias in
864 human reinforcement learning. *Brain Res* **1153**:111–121. doi:10.1016/j.brainres.2007.03.057
- 865 Bogacz R, Wagenmakers E-J, Forstmann BU, Nieuwenhuis S. 2010. The neural basis of the speed-accuracy
866 tradeoff. *Trends Neurosci* **33**:10–6. doi:10.1016/j.tins.2009.09.002
- 867 Boucher L, Palmeri TJ, Logan GD, Schall JD. 2007. Inhibitory Control in Mind and Brain : An Interactive Race

- 868 Model of Countermanding Saccades **114**:376–397. doi:10.1037/0033-295X.114.2.376
- 869 Brooks SP, Gelman A. 1998. General Methods for Monitoring Convergence of Iterative Simulations. *J Comput*
- 870 *Graph Stat* **7**:434–455. doi:10.1080/10618600.1998.10474787
- 871 Brown SD, Heathcote A. 2008. The simplest complete model of choice response time: Linear ballistic
- 872 accumulation. *Cogn Psychol* **57**:153–178. doi:10.1016/j.cogpsych.2007.12.002
- 873 Christakou A, Gershman SJ, Niv Y, Simmons A, Brammer M, Rubia K. 2013. Neural and Psychological
- 874 Maturation of Decision-making in Adolescence and Young Adulthood. *J Cogn Neurosci* **25**:1807–1823.
- 875 doi:10.1162/jocn_a_00447
- 876 Cisek P, Puskas GA, El-Murr S. 2009. Decisions in changing conditions: the urgency-gating model. *J Neurosci*
- 877 **29**:11560–71. doi:10.1523/JNEUROSCI.1844-09.2009
- 878 Collins AGE, Frank MJ. 2018. Within- and across-trial dynamics of human EEG reveal cooperative interplay
- 879 between reinforcement learning and working memory. *Proc Natl Acad Sci* **115**:2502–2507.
- 880 doi:10.1073/pnas.1720963115
- 881 Collins AGE, Frank MJ. 2012a. How much of reinforcement learning is working memory, not reinforcement
- 882 learning? A behavioral, computational, and neurogenetic analysis. *Eur J Neurosci* **35**:1024–1035.
- 883 doi:10.1111/j.1460-9568.2011.07980.x
- 884 Collins AGE, Frank MJ. 2012b. How much of reinforcement learning is working memory, not reinforcement
- 885 learning? A behavioral, computational, and neurogenetic analysis. *Eur J Neurosci* **35**:1024–1035.
- 886 doi:10.1111/j.1460-9568.2011.07980.x
- 887 Costa VD, Tran VL, Turchi J, Averbeck BB. 2015. Reversal learning and dopamine: A Bayesian perspective. *J*
- 888 *Neurosci* **35**:2407–2416. doi:10.1523/JNEUROSCI.1989-14.2015
- 889 Daw ND, Dayan P. 2014. The algorithmic anatomy of model-based evaluation. *Philos Trans R Soc B Biol Sci*
- 890 **369**. doi:10.1098/rstb.2013.0478
- 891 Daw ND, Kakade S, Dayan P. 2002. Opponent interactions between serotonin and dopamine. *Neural Networks*
- 892 **15**:603–616. doi:10.1016/S0893-6080(02)00052-7
- 893 Daw ND, O’Doherty JP, Dayan P, Seymour B, Dolan RJ. 2006. Cortical substrates for exploratory decisions in
- 894 humans. *Nature* **441**:876–879. doi:10.1038/nature04766
- 895 Dayan P, Daw ND. 2008. Decision theory, reinforcement learning, and the brain. *Cogn Affect Behav Neurosci*
- 896 **8**:429–453. doi:10.3758/CABN.8.4.429
- 897 Donkin C, Brown SD. 2018. Response Times and Decision-Making, Stevens’ Handbook of Experimental
- 898 Psychology and Cognitive Neuroscience. doi:10.1002/9781119170174.epcn509
- 899 Donkin C, Brown SD, Heathcote A. 2011. Drawing conclusions from choice response time models: A tutorial
- 900 using the linear ballistic accumulator. *J Math Psychol* **55**:140–151. doi:10.1016/j.jmp.2010.10.001
- 901 Donkin C, Brown SD, Heathcote A. 2009. The overconstraint of response time models: Rethinking the scaling
- 902 problem. *Psychon Bull Rev* **16**:1129–1135. doi:10.3758/PBR.16.6.1129
- 903 Dutilh G, Rieskamp J. 2016. Comparing perceptual and preferential decision making. *Psychon Bull Rev* **23**:723–
- 904 737. doi:10.3758/s13423-015-0941-1
- 905 Evans NJ, Brown SD, Mewhort DJK, Heathcote A. 2018. Refining the law of practice. *Psychol Rev* **125**:592–
- 906 605. doi:10.1037/rev0000105
- 907 Fontanesi L, Gluth S, Spektor MS, Rieskamp J. 2019a. A reinforcement learning diffusion decision model for
- 908 value-based decisions. *Psychon Bull Rev*. doi:10.3758/s13423-018-1554-2
- 909 Fontanesi L, Palminteri S, Lebreton M. 2019b. Decomposing the effects of context valence and feedback
- 910 information on speed and accuracy during reinforcement learning: a meta-analytical approach using
- 911 diffusion decision modeling. *Cogn Affect Behav Neurosci* **19**:490–502. doi:10.3758/s13415-019-00723-1
- 912 Forstmann BU, Ratcliff R, Wagenmakers E-J. 2016. Sequential Sampling Models in Cognitive Neuroscience:
- 913 Advantages, Applications, and Extensions. *Annu Rev Psychol* **67**:641–666. doi:10.1146/annurev-psych-
- 914 122414-033645
- 915 Frank MJ. 2004. By Carrot or by Stick: Cognitive Reinforcement Learning in Parkinsonism. *Science (80-)*
- 916 **306**:1940–1943. doi:10.1126/science.1102941
- 917 Frank MJ, Doll BB, Oas-Terpstra J, Moreno F. 2009. Prefrontal and striatal dopaminergic genes predict
- 918 individual differences in exploration and exploitation. *Nat Neurosci* **12**:1062–1068. doi:10.1038/nn.2342
- 919 Gelman A, Hill J. 2007. Data Analysis Using Regression and Multilevel/Hierarchical Models. Cambridge:
- 920 Cambridge University Press.
- 921 Gelman A, Rubin DB. 1992. Inference from Iterative Simulation Using Multiple Sequences. *Stat Sci* **7**:457–472.
- 922 doi:10.1214/ss/1177011136
- 923 Gershman SJ. 2015. Do learning rates adapt to the distribution of rewards? *Psychon Bull Rev* **22**:1320–1327.
- 924 doi:10.3758/s13423-014-0790-3
- 925 Haughey HM, Hutchison KE, Curran T, Frank MJ, Moustafa AA. 2007. Genetic triple dissociation reveals
- 926 multiple roles for dopamine in reinforcement learning. *Proc Natl Acad Sci* **104**:16311–16316.
- 927 doi:10.1073/pnas.0706111104

- 928 Hawkins GE, Forstmann BU, Wagenmakers E-J, Ratcliff R, Brown SD. 2015. Revisiting the Evidence for
929 Collapsing Boundaries and Urgency Signals in Perceptual Decision-Making. *J Neurosci* **35**:2476–2484.
930 doi:10.1523/JNEUROSCI.2410-14.2015
- 931 Hawkins GE, Heathcote A. 2020. Racing Against The Clock: Evidence-Based Vs. Time-Based Decisions.
932 *Psychol Rev*.
- 933 Heathcote A, Brown S, Mewhort DJK. 2000. The power law repealed: The case for an exponential law of
934 practice. *Psychon Bull Rev* **7**:185–207. doi:10.3758/BF03212979
- 935 Heathcote A, Brown SD, Wagenmakers E-J. 2015. An Introduction to Good Practices in Cognitive Modeling
936 Introduction to Model-Based Cognitive Neuroscience. New York, NY: Springer New York. pp. 25–48.
937 doi:10.1007/978-1-4939-2236-9_2
- 938 Heathcote A, Lin YS, Reynolds A, Strickland L, Gretton M, Matzke D. 2019. Dynamic models of choice. *Behav*
939 *Res Methods* **51**:961–985. doi:10.3758/s13428-018-1067-y
- 940 Heathcote A, Love J. 2012. Linear deterministic accumulator models of simple choice. *Front Psychol* **3**:1–19.
941 doi:10.3389/fpsyg.2012.00292
- 942 Ho T, Brown SD, Van Maanen L, Forstmann BU, Wagenmakers E-J, Serences JT. 2012. The Optimality of
943 Sensory Processing during the Speed-Accuracy Tradeoff. *J Neurosci* **32**:7992–8003.
944 doi:10.1523/JNEUROSCI.0340-12.2012
- 945 Izquierdo A, Brigman JL, Radke AK, Rudebeck PH, Holmes A. 2017. The neural basis of reversal learning: An
946 updated perspective. *Neuroscience* **345**:12–26. doi:10.1016/j.neuroscience.2016.03.021
- 947 Jang AI, Costa VD, Rudebeck PH, Chudasama Y, Murray EA, Averbeck BB. 2015. The role of frontal cortical
948 and medial-temporal lobe brain areas in learning a Bayesian prior belief on reversals. *J Neurosci*
949 **35**:11751–11760. doi:10.1523/JNEUROSCI.1594-15.2015
- 950 Katsimpokis D, Hawkins GE, Van Maanen L. 2020. Not all Speed-Accuracy Trade-Off Manipulations Have the
951 Same Psychological Effect. *Comput Brain Behav*. doi:10.1007/s42113-020-00074-y
- 952 Kuznetsova A, Brockhoff PB, Christensen RHB. 2017. lmerTest Package: Tests in Linear Mixed Effects
953 Models. *J Stat Softw* **82**. doi:10.18637/jss.v082.i13
- 954 Leite FP, Ratcliff R. 2010. Modeling reaction time and accuracy of multiple-alternative decisions. *Attention,*
955 *Perception, Psychophys* **72**:246–273. doi:10.3758/APP.72.1.246
- 956 Logan GD, Van Zandt T, Verbruggen F, Wagenmakers EJ. 2014. On the ability to inhibit thought and action:
957 General and special theories of an act of control. *Psychol Rev* **121**:66–95. doi:10.1037/a0035230
- 958 Luzardo A, Alonso E, Mondragón E. 2017. A Rescorla-Wagner drift-diffusion model of conditioning and
959 timing. *PLOS Comput Biol* **13**:e1005796. doi:10.1371/journal.pcbi.1005796
- 960 McDougle SD, Collins AGE. 2020. Modeling the influence of working memory, reinforcement, and action
961 uncertainty on reaction time and choice during instrumental learning. *Psychon Bull Rev*.
962 doi:10.3758/s13423-020-01774-z
- 963 Miletić S. 2016. Neural Evidence for a Role of Urgency in the Speed-Accuracy Trade-off in Perceptual
964 Decision-Making. *J Neurosci* **36**:5909–5910. doi:10.1523/JNEUROSCI.0894-16.2016
- 965 Miletić S, Boag RJ, Forstmann BU. 2020. Mutual benefits: Combining reinforcement learning with sequential
966 sampling models. *Neuropsychologia* **136**. doi:10.1016/j.neuropsychologia.2019.107261
- 967 Miletić S, Turner BM, Forstmann BU, Van Maanen L. 2017. Parameter recovery for the Leaky Competing
968 Accumulator model. *J Math Psychol* **76**:25–50. doi:10.1016/j.jmp.2016.12.001
- 969 Miletić S, Van Maanen L. 2019. Caution in decision-making under time pressure is mediated by timing ability.
970 *Cogn Psychol* **110**:16–29. doi:10.1016/j.cogpsych.2019.01.002
- 971 Millner AJ, Gershman SJ, Nock MK, den Ouden HEM. 2018. Pavlovian Control of Escape and Avoidance. *J*
972 *Cogn Neurosci* **30**:1379–1390. doi:10.1162/jocn_a_01224
- 973 Moran R. 2016. Thou shalt identify! The identifiability of two high-threshold models in confidence-rating
974 recognition (and super-recognition) paradigms. *J Math Psychol* **73**:1–11. doi:10.1016/j.jmp.2016.03.002
- 975 Murphy PR, Boonstra E, Nieuwenhuis S. 2016. Global gain modulation generates time-dependent urgency
976 during perceptual choice in humans. *Nat Commun* **7**:1–14. doi:10.1038/ncomms13526
- 977 Niv Y, Edlund JA, Dayan P, O’Doherty JP. 2012. Neural Prediction Errors Reveal a Risk-Sensitive
978 Reinforcement-Learning Process in the Human Brain. *J Neurosci* **32**:551–562.
979 doi:10.1523/jneurosci.5498-10.2012
- 980 O’Doherty JP, Cockburn J, Pauli WM. 2017. Learning, Reward, and Decision Making. *Annu Rev Psychol*
981 **68**:73–100. doi:10.1146/annurev-psych-010416-044216
- 982 Pachella RG, Pew RW. 1968. Speed-Accuracy Tradeoff in Reaction Time: Effect of Discrete Criterion Times. *J*
983 *Exp Psychol* **76**:19–24. doi:10.1037/h0021275
- 984 Palminteri S, Khamassi M, Joffily M, Coricelli G. 2015. Contextual modulation of value signals in reward and
985 punishment learning. *Nat Commun* **6**. doi:10.1038/ncomms9096
- 986 Palminteri S, Wyart V, Koehlin E. 2017. The Importance of Falsification in Computational Cognitive
987 Modeling. *Trends Cogn Sci* **21**:425–433. doi:10.1016/j.tics.2017.03.011

- 988 Pedersen ML, Frank MJ. 2020. Simultaneous Hierarchical Bayesian Parameter Estimation for Reinforcement
989 Learning and Drift Diffusion Models: a Tutorial and Links to Neural Data. *Comput Brain Behav*.
990 doi:10.1007/s42113-020-00084-w
- 991 Pedersen ML, Frank MJ, Biele G. 2017. The drift diffusion model as the choice rule in reinforcement learning.
992 *Psychon Bull Rev* **24**:1234–1251. doi:10.3758/s13423-016-1199-y
- 993 Peirce J, Gray JR, Simpson S, MacAskill M, Höchenberger R, Sogo H, Kastman E, Lindeløv JK. 2019.
994 PsychoPy2: Experiments in behavior made easy. *Behav Res Methods* **51**:195–203. doi:10.3758/s13428-
995 018-01193-y
- 996 Purcell BA, Heitz RP, Cohen JY, Schall JD, Logan GD, Palmeri TJ. 2010. Neurally constrained modeling of
997 perceptual decision making. *Psychol Rev* **117**:1113–1143. doi:10.1037/a0020311
- 998 R Core Team. 2017. R: A language and environment for statistical computing.
- 999 Rae B, Heathcote A, Donkin C, Averell L, Brown SD. 2014. The hare and the tortoise: Emphasizing speed can
1000 change the evidence used to make decisions. *J Exp Psychol Learn Mem Cogn* **1**–39.
1001 doi:10.1037/a0036801
- 1002 Ratcliff R. 1978. A theory of memory retrieval. *Psychol Rev* **85**:59–108.
- 1003 Ratcliff R, Hasegawa YT, Hasegawa RP, Childers R, Smith PL, Segraves MA. 2011. Inhibition in superior
1004 colliculus neurons in a brightness discrimination task? *Neural Comput* **23**:1790–1820.
1005 doi:10.1162/NECO_a_00135
- 1006 Ratcliff R, Hasegawa YT, Hasegawa RP, Smith PL, Segraves MA. 2007. Dual Diffusion Model for Single-Cell
1007 Recording Data From the Superior Colliculus in a Brightness-Discrimination Task. *J Neurophysiol*
1008 **97**:1756–1774. doi:10.1152/jn.00393.2006
- 1009 Ratcliff R, McKoon G. 2008. The diffusion decision model: theory and data for two-choice decision tasks.
1010 *Neural Comput* **20**:873–922. doi:10.1162/neco.2008.12-06-420
- 1011 Ratcliff R, Rouder JN. 1998. Modeling Response Times for Two-Choice Decisions. *Psychol Sci* **9**:347–356.
- 1012 Ratcliff R, Smith PL, Brown SD, McKoon G. 2016. Diffusion Decision Model: Current Issues and History.
1013 *Trends Cogn Sci* **20**:260–281. doi:10.1016/j.tics.2016.01.007
- 1014 Ratcliff R, Voskuilen C, Teodorescu A. 2018. Modeling 2-alternative forced-choice tasks: Accounting for both
1015 magnitude and difference effects. *Cogn Psychol* **103**:1–22. doi:10.1016/j.cogpsych.2018.02.002
- 1016 Rescorla RA, Wagner AR. 1972. A theory of Pavlovian conditioning: Variations in the effectiveness of
1017 reinforcement and nonreinforcement. *Class Cond II Curr Res Theory* **21**:64–99.
1018 doi:10.1101/gr.110528.110
- 1019 Rummery GA, Niranjan M. 1994. On-Line Q-Learning Using Connectionist Systems.
- 1020 Satterthwaite FE. 1941. Synthesis of variance. *Psychometrika* **6**:309–316. doi:10.1007/BF02288586
- 1021 Sewell DK, Jach HK, Boag RJ, Van Heer CA. 2019. Combining error-driven models of associative learning
1022 with evidence accumulation models of decision-making. *Psychon Bull Rev*. doi:10.3758/s13423-019-
1023 01570-4
- 1024 Sewell DK, Stallman A. 2020. Modeling the Effect of Speed Emphasis in Probabilistic Category Learning.
1025 *Comput Brain Behav* **3**:129–152. doi:10.1007/s42113-019-00067-6
- 1026 Shahar N, Hauser TU, Moutoussis M, Moran R, Keramati M, Consortium NSPN, Dolan RJ. 2019. Improving
1027 the reliability of model-based decision-making estimates in the two-stage decision task with reaction-
1028 times and drift-diffusion modeling. *PLoS Comput Biol* **15**:1–25. doi:10.1371/journal.pcbi.1006803
- 1029 Spektor MS, Kellen D. 2018. The relative merit of empirical priors in non-identifiable and sloppy models:
1030 Applications to models of learning and decision-making. *Psychon Bull Rev* **25**:2047–2068.
1031 doi:10.3758/s13423-018-1446-5
- 1032 Spiegelhalter DJ, Best NG, Carlin BP, Van der Linde A. 2002. Bayesian measures of model complexity and fit.
1033 *J R Stat Soc Ser B (Statistical Methodol)* **64**:583–639.
- 1034 Sutton, Richard S. 1988. Learning to Predict by the Method of Temporal Differences. *Mach Learn* **3**:9–44.
1035 doi:10.1023/A:1018056104778
- 1036 Sutton RS, Barto AG. 2018. Reinforcement Learning: An Introduction, 2nd ed, MIT Press. Cambridge, MA:
1037 MIT press.
- 1038 Teodorescu AR, Moran R, Usher M. 2016. Absolutely relative or relatively absolute: violations of value
1039 invariance in human decision making. *Psychon Bull Rev* **23**:22–38. doi:10.3758/s13423-015-0858-8
- 1040 Ter Braak CJF. 2006. A Markov Chain Monte Carlo version of the genetic algorithm Differential Evolution:
1041 easy Bayesian computing for real parameter spaces. *Stat Comput* **16**:239–249. doi:10.1007/s11222-006-
1042 8769-1
- 1043 Thura D, Cisek P. 2016. Modulation of Premotor and Primary Motor Cortical Activity during Volitional
1044 Adjustments of Speed-Accuracy Trade-Offs. *J Neurosci* **36**:938–956. doi:10.1523/JNEUROSCI.2230-
1045 15.2016
- 1046 Tillman G, Van Zandt T, Logan GD. 2020. Sequential sampling models without random between-trial
1047 variability: the racing diffusion model of speeded decision making. *Psychon Bull Rev*.

- 1048 doi:10.3758/s13423-020-01719-6
- 1049 Tran H, Van Maanen L, Matzke D, Heathcote A. n.d. Systematic parameter reviews in cognitive modeling:
1050 Towards robust and cumulative models of psychological processes.
- 1051 Trueblood JS, Heathcote A, Evans NJ, Holmes WR. 2020. Urgency, Leakage, and the Relative Nature of
1052 Information Processing in Decision-making. *Psychol Rev* 706291. doi:10.1101/706291
- 1053 Turner BM. 2019. Toward a Common Representational Framework for Adaptation. *Psychol Rev*.
1054 doi:10.1037/rev0000148
- 1055 Turner BM, Sederberg PB, Brown SD, Steyvers M. 2013. A method for efficiently sampling from distributions
1056 with correlated dimensions. *Psychol Methods* **18**:368–384. doi:10.1037/a0032222
- 1057 van Maanen L, Miletic S. 2020. The interpretation of behavior-model correlations in unidentified cognitive
1058 models. *Psychon Bull Rev*. doi:10.3758/s13423-020-01783-y
- 1059 van Maanen L, van der Mijl R, van Beurden MHPH, Roijendijk LMM, Kingma BRM, Miletic S, van Rijn H.
1060 2019. Core body temperature speeds up temporal processing and choice behavior under deadlines. *Sci Rep*
1061 **9**:10053. doi:10.1038/s41598-019-46073-3
- 1062 Van Ravenzwaaij D, Brown SD, Marley AAJ, Heathcote A. 2020. Accumulating advantages: A new
1063 conceptualization of rapid multiple choice. *Psychol Rev* **127**:186–215. doi:10.1037/rev0000166
- 1064 Voss A, Nagler M, Lerche V. 2013. Diffusion models in experimental psychology: A practical introduction. *Exp*
1065 *Psychol* **60**:385–402. doi:10.1027/1618-3169/a000218
- 1066 Voss A, Rothermund K, Voss J. 2004. Interpreting the parameters of the diffusion model: An empirical
1067 validation. *Mem Cognit* **32**:1206–1220. doi:10.3758/BF03196893
- 1068
- 1069
- 1070

1071 **Figure supplements**

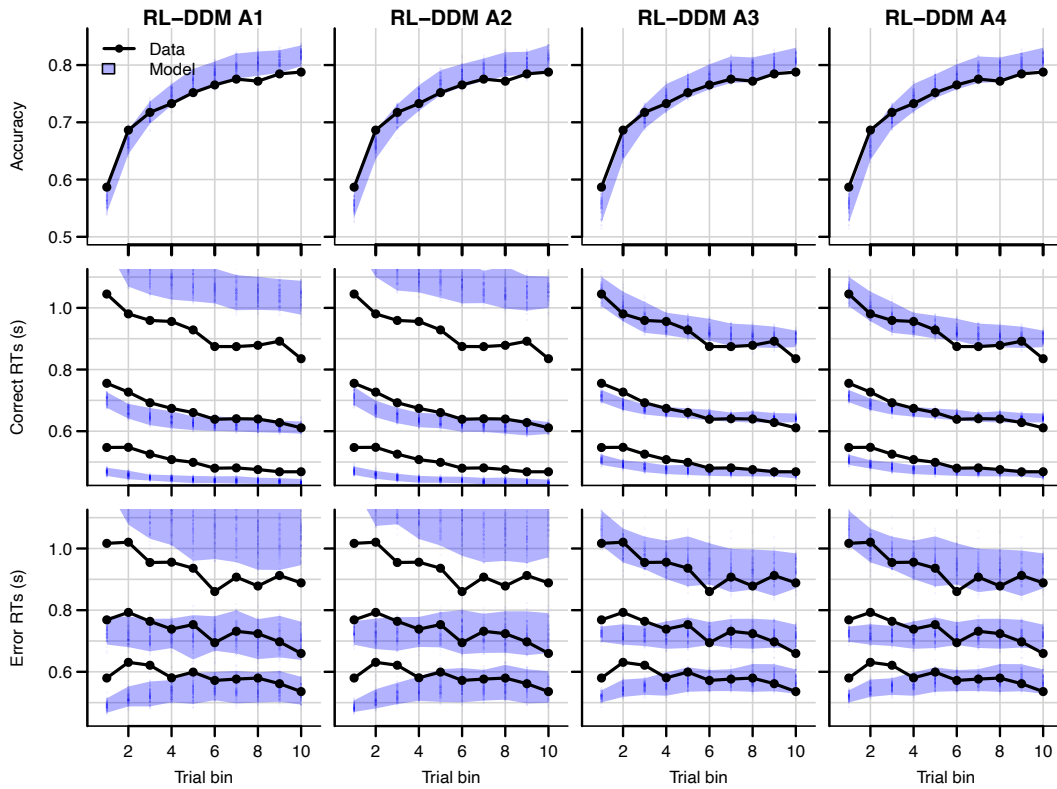


Figure 3-figure supplement 1. Comparison of posterior predictive distributions of four additional RL-DDMs. Data are black dots and lines, posterior predictive distribution are blue. Top row depicts accuracy over trial bins. Middle and bottom row illustrate 10th, 50th and 90th quantile RT for the correct (middle row) and error (bottom row) response over trial bins. Shaded areas correspond to the 95% credible interval of the posterior predictive distributions. All data are collapsed across participants and difficulty conditions. The summed BPICs were 7717 (RL-DDM A1), 7636 (RL-DDM A2), 4844 (RL-DDM A3) and 4884 (RL-DDM A4). Hence, the largest improvement of quality of fit of the RL-DDM was obtained by adding s_{t0} .

1072

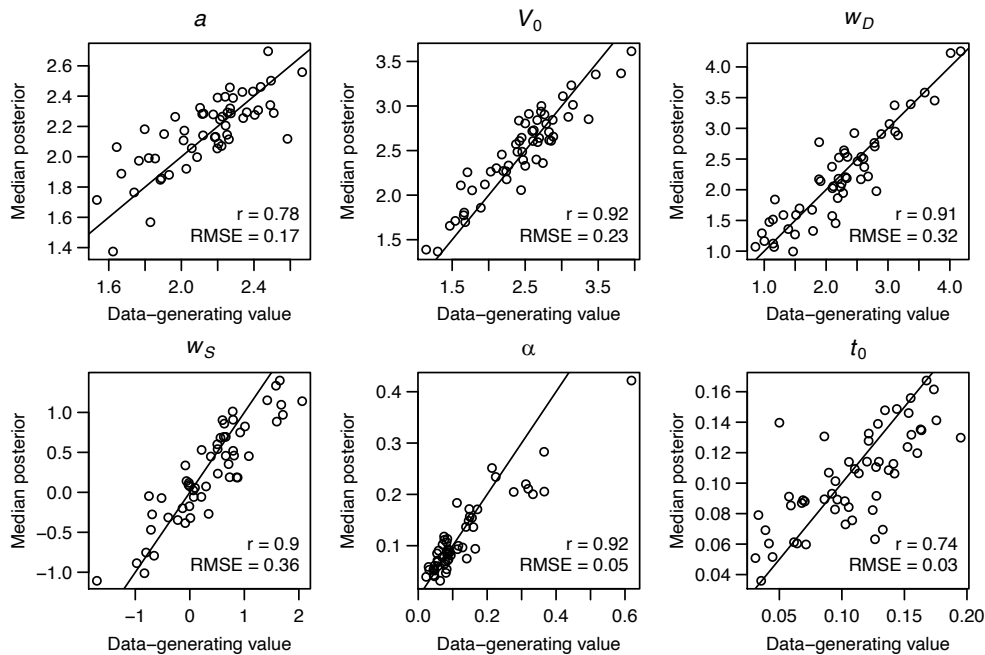


Figure 3-figure supplement 2. Parameter recovery of the RL-ARD model, using the experimental paradigm of experiment 1. Parameter recovery was done by first fitting the RL-ARD model to the empirical data, and then simulating the exact same experimental paradigm (208 trials, 55 subjects, 4 difficulty conditions) using the median parameter estimates obtained from the model fit. Subsequently, the RL-ARD was fit to the simulated data. The recovered median posterior estimates (y-axis) are plotted against the data-generating values (x-axis). Pearson's correlation coefficient r and the root mean square error (RMSE) are shown in each panel. Diagonal lines indicate the identity $x = y$.

1073
1074

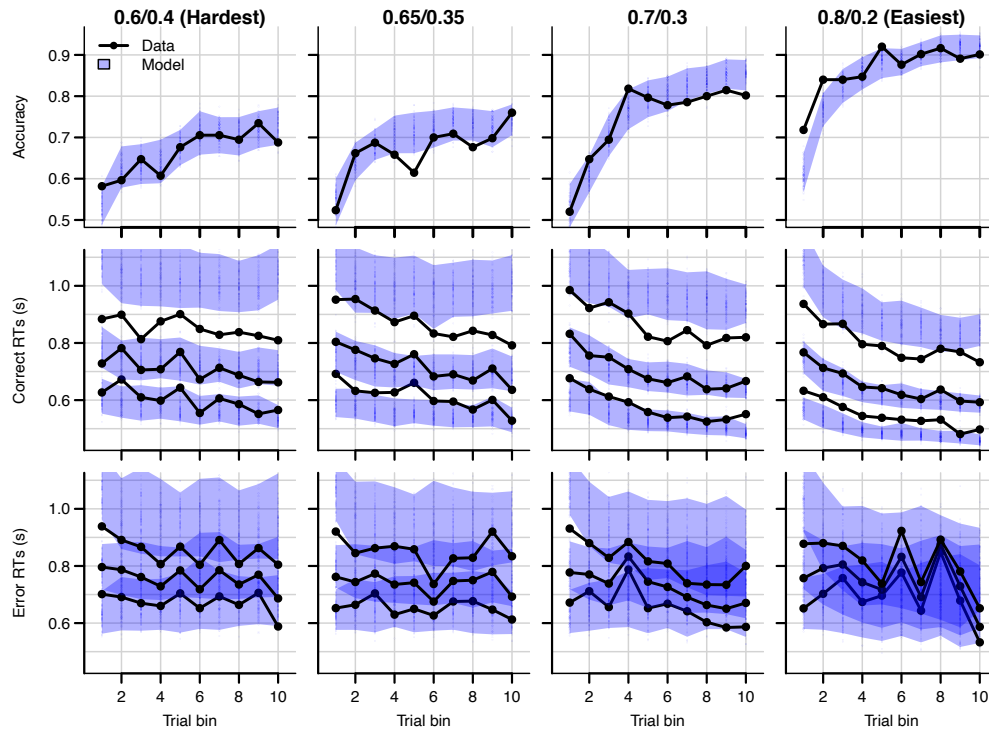


Figure 4-figure supplement 1. Data (black) and posterior predictive distribution of the RL-DDM (blue), separately for each difficulty condition. Row titles indicate the reward probabilities, with 0.6/0.4 being the most difficult, and 0.8/0.2 the easiest condition. Top row depicts accuracy over trial bins. Middle and bottom row illustrate 10th, 50th, and 90th quantile RT for the correct (middle row) and error (bottom row) response over trial bins. Shaded areas correspond to the 95% credible interval of the posterior predictive distributions. All data are collapsed across participants.

1075
1076

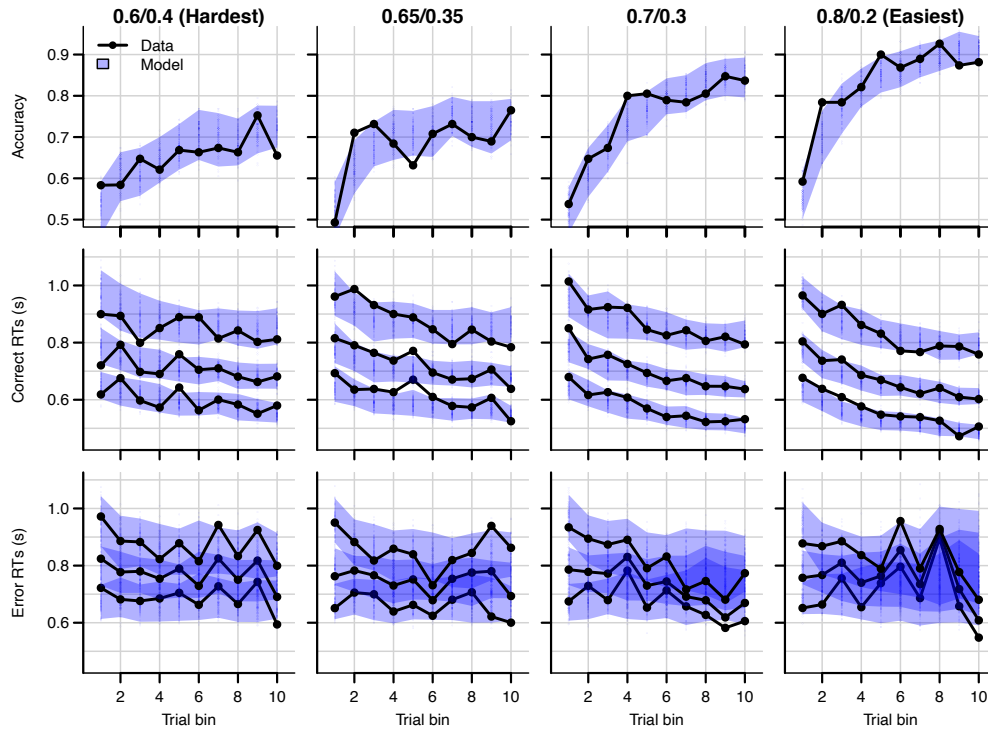


Figure 4-figure supplement 2. Data (black) and posterior predictive distribution of the RL-ARD (blue), separately for each difficulty condition, excluding 17 subjects which had perfect accuracy in the first bin of the easiest condition. Row titles indicate the reward probabilities, with 0.6/0.4 being the most difficult, and 0.8/0.2 the easiest condition. Top row depicts accuracy over trial bins. Middle and bottom row illustrate 10th, 50th, and 90th quantile RT for the correct (middle row) and error (bottom row) response over trial bins. Shaded areas correspond to the 95% credible interval of the posterior predictive distributions. All data are collapsed across participants.

1077
1078
1079

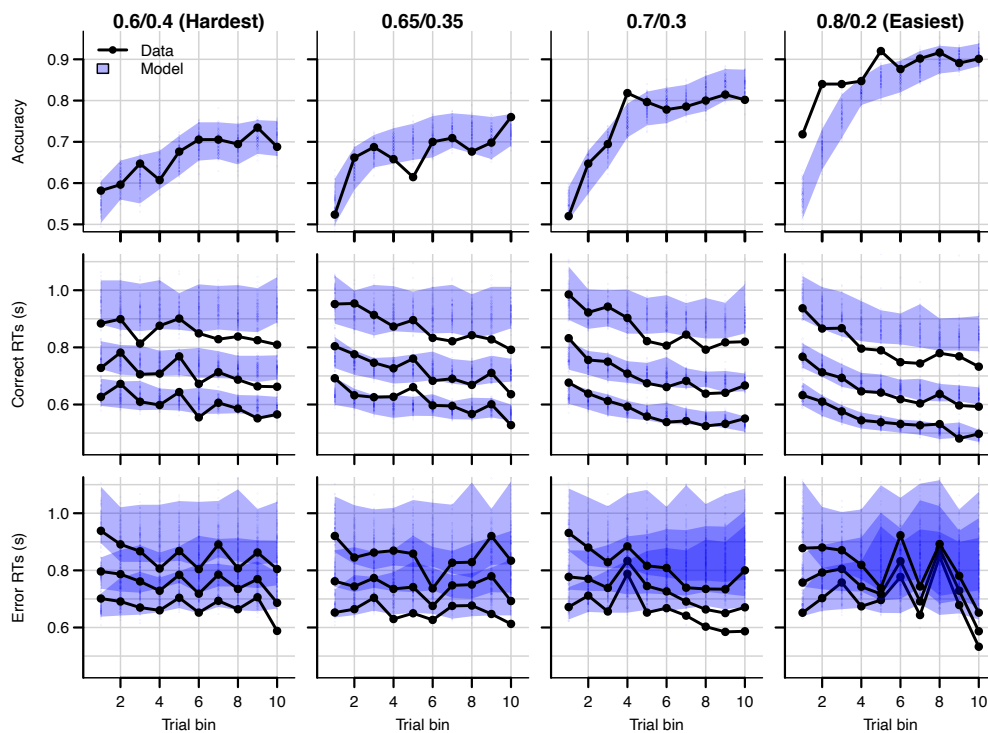


Figure 4-figure supplement 3. Posterior predictive distribution of the RL-ALBA model on the data of experiment 1, with one column per difficulty condition. The LBA assumes that, on every trial, two accumulators race deterministically towards a common bound b . Each accumulator i starts at a start point sampled from a uniform distribution $[0, A]$, and with a speed of evidence accumulation sampled from a normal distribution $\mathcal{N}(v_i, s_i)$. In the RL-ALBA model, we used Equation 3 to link Q-values to LBA drift rates v_1 and v_2 (excluding the sW term, since the LBA assumes no within-trial noise). Instead of directly estimating threshold b , we estimated the difference $B = b - A$ (which simplifies enforcing $b > A$). We used the following mildly informed priors for the hypermeans: $V_0 \sim \mathcal{N}(2, 5)$, $w_d \sim \mathcal{N}(9, 5)$ truncated to lower bound 0, $w_s \sim \mathcal{N}(0, 3)$, $s_2 \sim \mathcal{N}(1, 1)$, $A \sim \mathcal{N}(1, 1)$, $B \sim \mathcal{N}(3, 5)$ truncated to lower bound 0, and $t_0 \sim \mathcal{N}(0.3, 0.5)$ truncated to lower bound 0.025 and upper bound 1. For the hyperSDs, all priors were $\Gamma(1, 1)$. The summed BPIC was 4836, indicating that the RL-ALBA performs slightly better than the RL-DDM with between-trial variabilities (BPIC = 4844), and better than the RL-IARD (BPIC = 4849), but not as well as the RL-ARD (BPIC = 4577).

1080
1081
1082

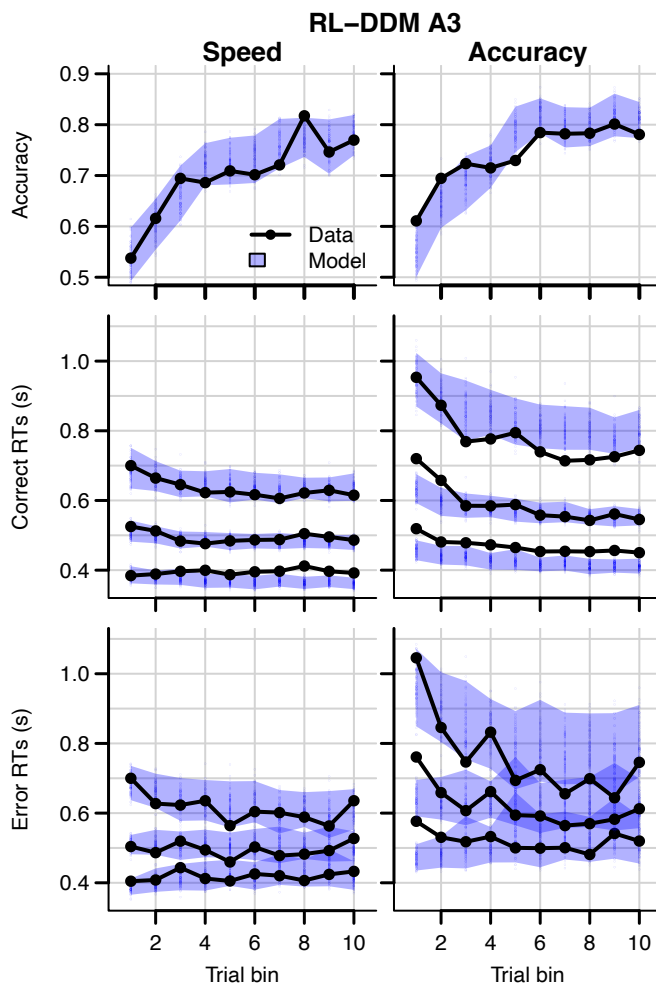


Figure 6-figure supplement 1. Data (black) of experiment 2 and posterior predictive distribution (blue) of the RL-DDM A3 with separate thresholds for the SAT conditions, and between-trial variabilities in drift rates, start points, and non-decision times. The corresponding summed BPIC was -861, an improvement over the RL-DDM, but outperformed by the RL-ARD ($\Delta BPIC = 232$ in favor of the RL-ARD). Top row depicts accuracy over trial bins. Middle and bottom row illustrate 10th, 50th and 90th quantile RT for the correct (middle row) and error (bottom row) response over trial bins. Left and right column are speed and accuracy emphasis condition, respectively. Shaded areas correspond to the 95% credible interval of the posterior predictive distributions.

1083
1084

1085
1086

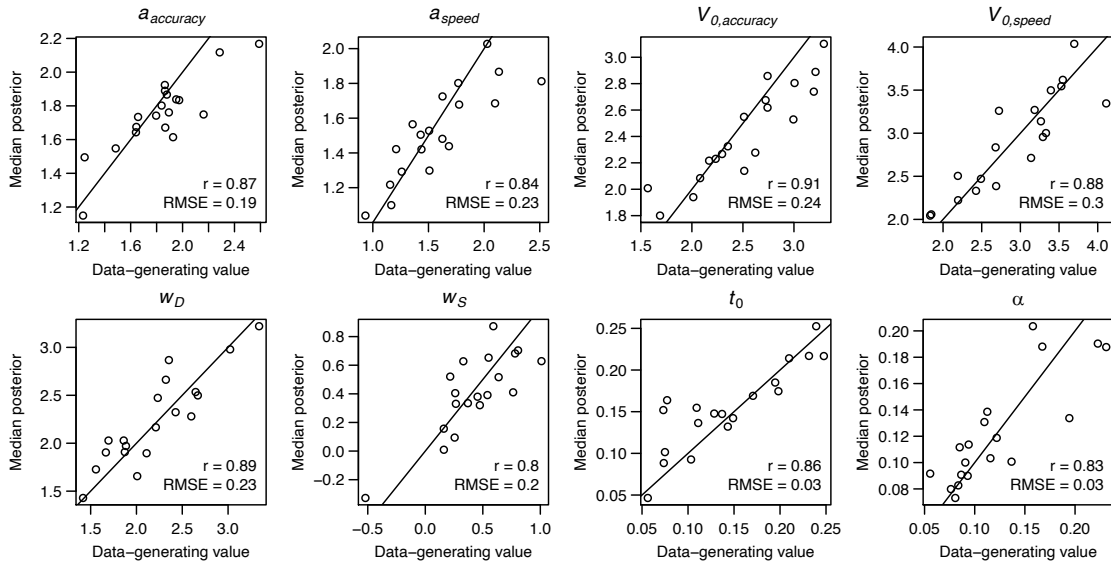


Figure 6-figure supplement 2. Parameter recovery of the RL-ARD model, using the experimental paradigm of experiment 2. Parameter recovery was done by first fitting the RL-ARD model to the empirical data, and then simulating the exact same experimental paradigm (19 subjects, 3 difficulty conditions, 2 SAT conditions, 312 trials) using the median parameter estimates obtained from the model fit. Subsequently, the RL-ARD was fit to the simulated data. The median posterior estimates (y-axis) are plotted against the data-generating values (x-axis). Pearson's correlation coefficient r and the root mean square error (RMSE) are shown in each panel. Diagonal lines indicate the identity $x = y$.

1087
1088
1089

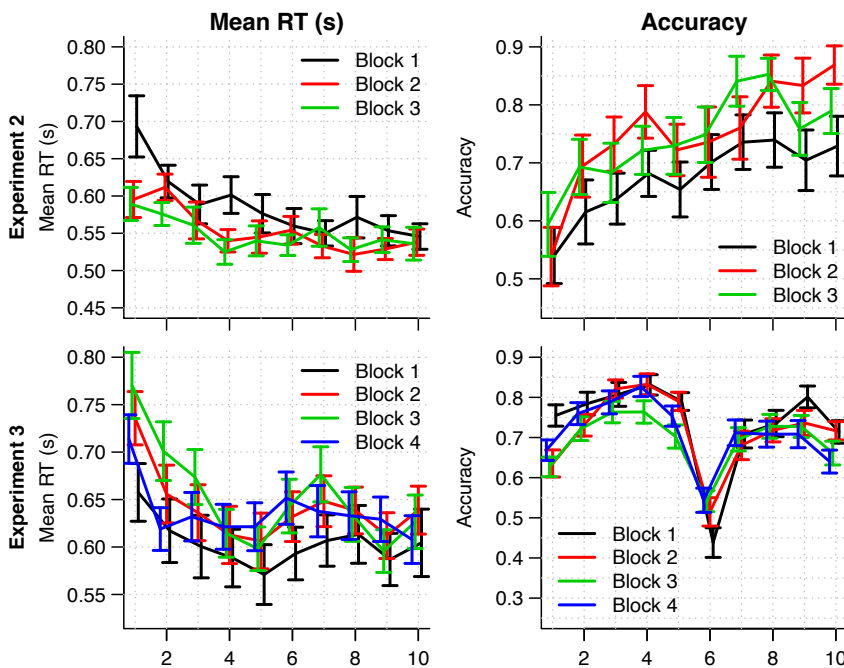


Figure 6-figure supplement 3. Mean RT (left column) and choice accuracy (right column) across trial bins (x-axis) for experiments 2 and 3 (rows). Block numbers are color-coded. Error bars are 1 SE.

Mixed effects models indicated that in experiment 2, RTs decreased with block number ($b = -0.04$, $SE = 6.15 \cdot 10^{-3}$, 95% CI $[-0.05, -0.03]$, $p = 6.61 \cdot 10^{-10}$) as well as with trial bin ($b = -0.02$, $SE = 2.11 \cdot 10^{-3}$, 95% CI $[-$

0.02, -0.01], $p = 1.68 \times 10^{-13}$), and there was an interaction between trial bin and block number ($\beta = 3.61 \times 10^{-3}$, $SE = 9.86 \times 10^{-4}$, 95% CI [0.00, 0.01], $p = 2.52 \times 10^{-4}$). There was a main effect of (log-transformed) trial bin on accuracy ($b = 0.36$, $SE = 0.11$, 95% CI [0.15, 0.57], $p = 7.99 \times 10^{-4}$), but no effect of block number, nor an interaction between block number and trial bin on accuracy.

In experiment 3, response times increased with block number ($b = 0.02$, $SE = 3.10 \times 10^{-3}$, 95% CI [0.01, 0.02], $p = 1.21 \times 10^{-7}$), decreased with trial bin ($b = -4.24 \times 10^{-3}$, $SE = 1.3 \times 10^{-3}$, 95% CI $[-6.92 \times 10^{-3}, -1.56 \times 10^{-3}]$, $p = 0.002$), but there was no interaction between trial bin and block number ($b = -9.15 \times 10^{-4}$, $SE = 5 \times 10^{-4}$, 95% CI [0.00, 0.00], $p = 0.067$). The bottom left panel suggests that the main effect of block number on RT is largely caused by an increase in RT after the first block. Accuracy decreased with (log-transformed) trial bin ($b = -0.12$, $SE = 0.05$, 95% CI $[-0.22, -0.02]$, $p = 0.02$), decreased with block number ($b = -0.08$, $SE = 0.03$, 95% CI $[-0.14, -0.02]$, $p = 0.009$), but there was no interaction ($b = 0.02$, $SE = 0.02$, 95% CI $[-0.02, 0.06]$, $p = 0.276$). The decrease in accuracy with trial bin is expected due to the presence of reversals. The combination of an increase in RT and a decrease in accuracy after the first block could indicate that participants learnt the structure of the task (i.e., the presence of reversals) in the first block, and adjusted their behavior accordingly. In line with this speculation, the accuracy in trial bin 6 (in which the reversal occurred) was lowest in the first block, which suggests that participants adjusted to the reversal faster in the later blocks.

1090
1091
1092

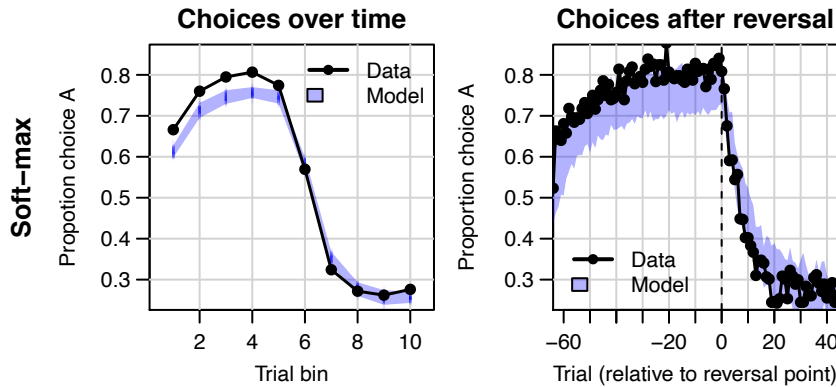


Figure 7-figure supplement 1. Data (black) of experiment 3 and posterior predictive of a standard soft-max learning model (blue). As priors, we used $\beta \sim N(1,5)$ truncated at 0 for the hypermean and $\Gamma(1,1)$ for the hyperSD. Left panel depicts choice proportions for option over trial bins, where choice A is defined as the high-probability reward choice prior to the reversal. Right column depicts choice proportion over trials, aligned to the trial at which the reversal occurred (trial 0). Shaded areas correspond to the 95% credible interval of the posterior predictive distributions.

1093
1094
1095
1096

RL-DDM A3

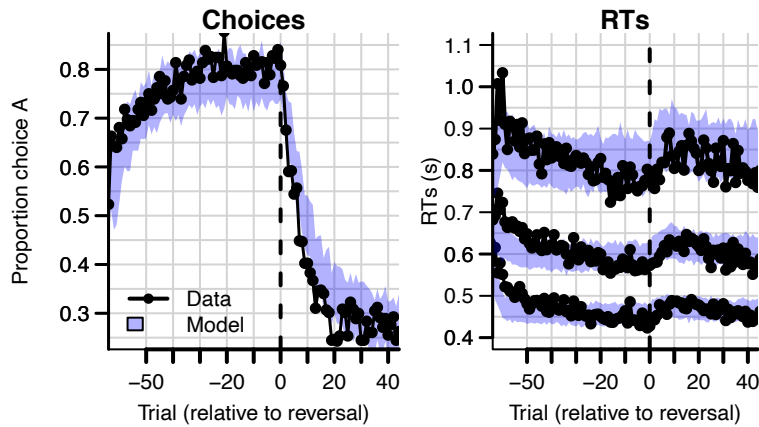


Figure 7-figure supplement 2. Data (black) of experiment 3 and posterior predictive distribution (blue) of the RL-DDM A3 (with between-trial variabilities in drift rates, start points, and non-decision times). The summed BPIC was 11659. This is better compared to the RL-DDM ($\Delta BPIC = 3940$) but did not outperform the RL-ARD ($\Delta BPIC = 112$ in favor of the RL-ARD).

1097

1098

1099

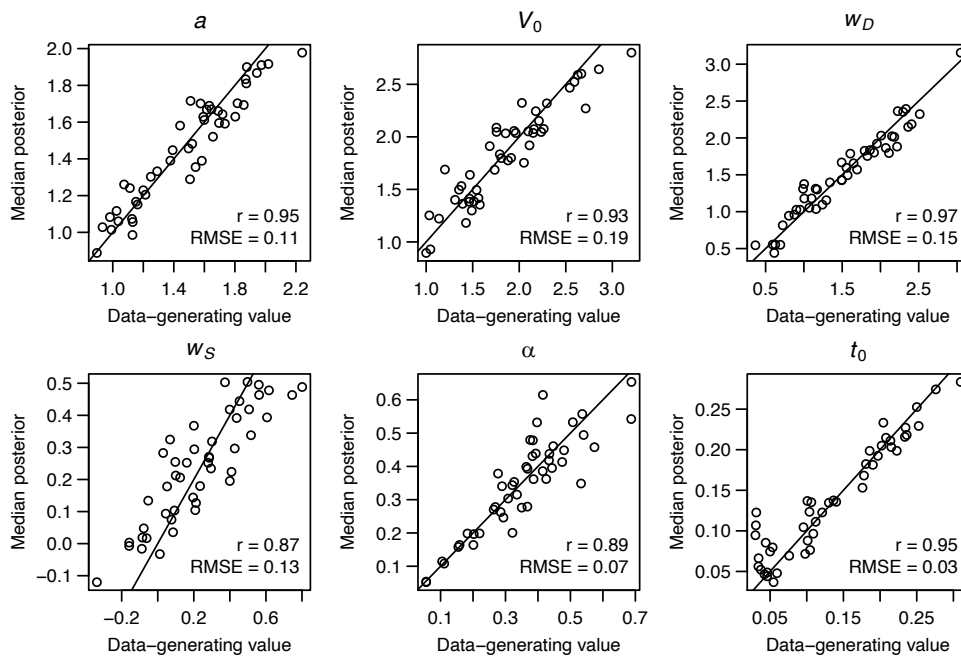


Figure 7-figure supplement 3. Parameter recovery of the RL-ARD model, using the experimental paradigm of experiment 3. Parameter recovery was done by first fitting the RL-ARD model to the empirical data, and then simulating the exact same experimental paradigm (49 subjects, 2 difficulty conditions, 512 trials including reversals) using the median parameter estimates obtained from the model fit. Subsequently, the RL-ARD was fit to the simulated data. The median posterior estimates (y-axis) are plotted against the data-generating values (x-axis). Pearson's correlation coefficient r and the root mean square error (RMSE) are shown in each panel. Diagonal lines indicate the identity $x = y$.

1100

1101

1102