# A new model of decision processing in instrumental learning tasks

Steven Miletić[a]*, Russell J. Boag[a], Anne C. Trutti[a,b], Birte U. Forstmann[a], Andrew Heathcote[a,c]

[a] University of Amsterdam, Department of Psychology, Nieuwe Achtergracht 129B, Amsterdam, The Netherlands
[b] Leiden University, Department of Psychology, Wassenaarseweg 52, Leiden, The Netherlands
[c] University of Newcastle, School of Psychology, Newcastle, Australia

*Correspondence concerning this article should be addressed to Steven Miletić, Nieuwe Achtergracht 129B, 1001NK Amsterdam, The Netherlands. E: s.miletic@uva.nl

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

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; Luzardo 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 Niranjan, 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; Miletić 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-lARD, 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
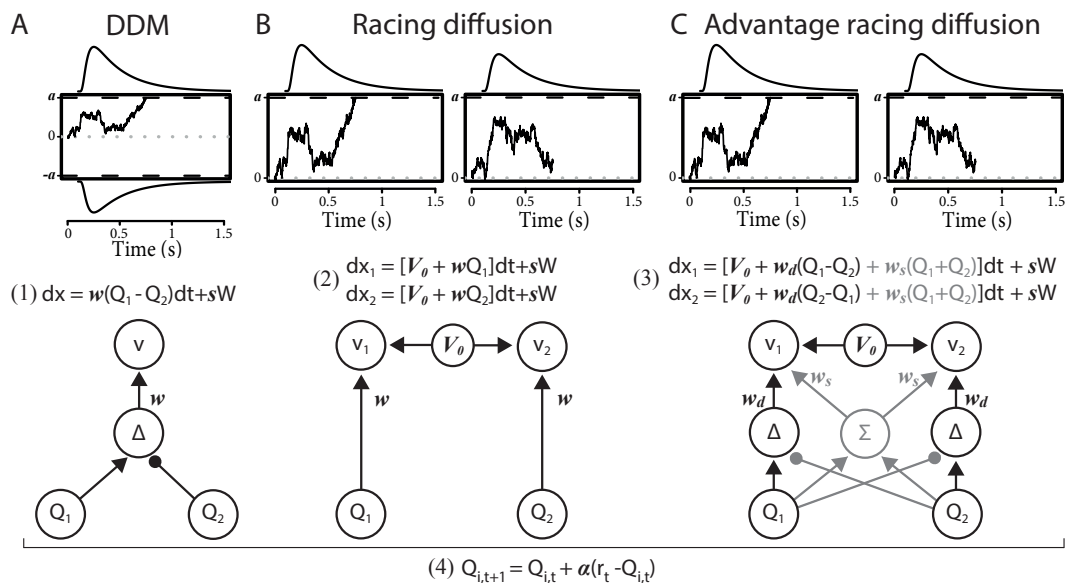139    distribution of RT, and learning-related changes thereof.



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 ($\Delta$) 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 ($\Delta$) 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* ($\Sigma$) 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 $\alpha$.

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 $10^{th}$, $50^{th}$ (median) and $90^{th}$
150     percentiles separately for correct and error responses. The median summarizes central
151     tendency, the difference between $10^{th}$ and $90^{th}$ percentiles summarizes variability and the larger
152     difference between the $90^{th}$ and $50^{th}$ percentiles than between the $50^{th}$ and $10^{th}$ 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),
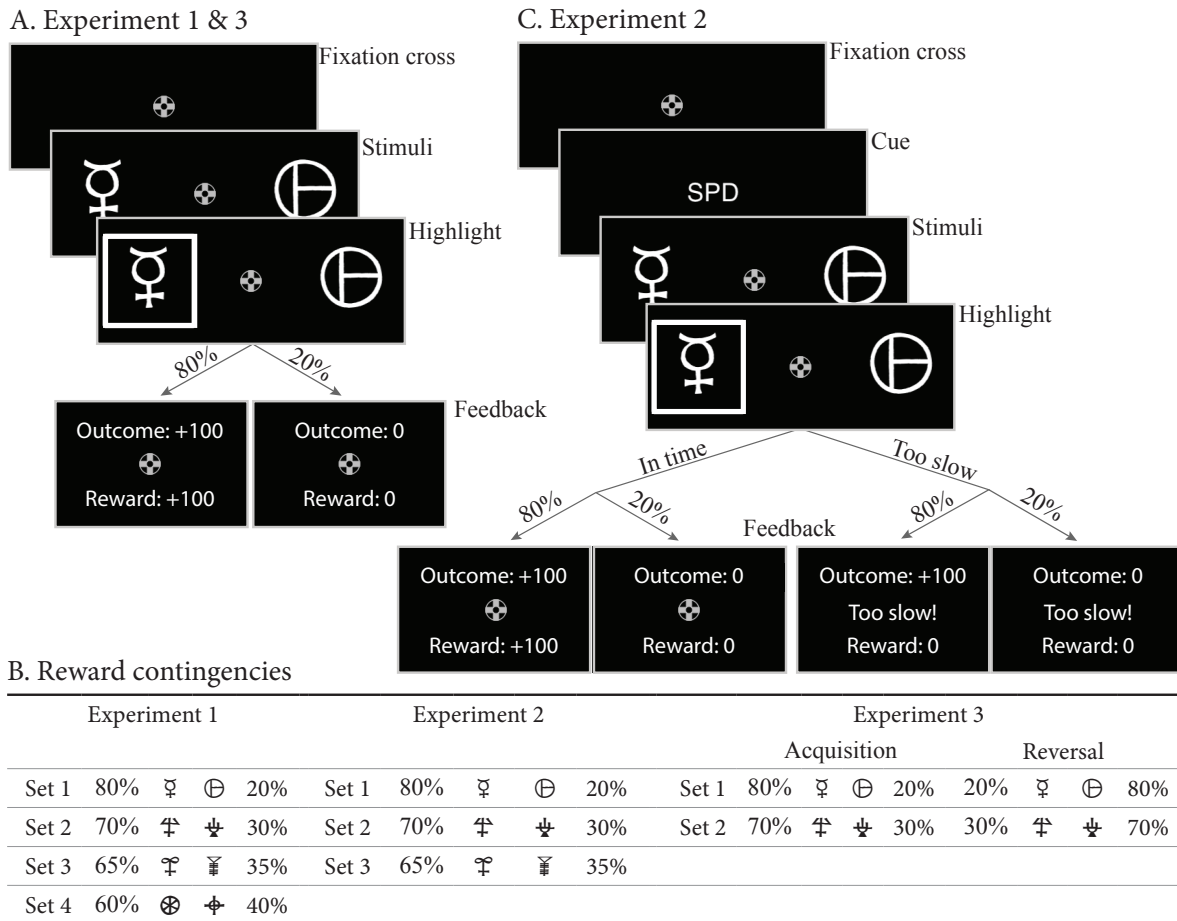163     ensuring that a better fit is not only due to greater model flexibility.
164

**A. Experiment 1 & 3**

Fixation cross

Stimuli

Highlight

80%  20%

Outcome: +100 / Reward: +100

Outcome: 0 / Reward: 0

Feedback

**C. Experiment 2**

Fixation cross

Cue

SPD

Stimuli

Highlight

In time     Too slow

80%  20%     80%  20%

| Outcome: +100 Reward: +100 | Outcome: 0 Reward: 0 | Outcome: +100 Too slow! Reward: 0 | Outcome: 0 Too slow! Reward: 0 |

Feedback

**B. Reward contingencies**

| | Experiment 1 | | | | Experiment 2 | | | | Experiment 3 | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | | | | | | | Acquisition | | | Reversal | |
| Set 1 | 80% | ☿ | ⊕ | 20% | Set 1 | 80% | ☿ | ⊕ | 20% | Set 1 | 80% | ☿ ⊕ | 20% | 20% ☿ ⊕ 80% |
| Set 2 | 70% | ♃ | ♇ | 30% | Set 2 | 70% | ♃ | ♇ | 30% | Set 2 | 70% | ♃ ♇ | 30% | 30% ♃ ♇ 70% |
| Set 3 | 65% | ♈ | ⚷ | 35% | Set 3 | 65% | ♈ | ⚷ | 35% | | | | | |
| Set 4 | 60% | ⊕ | ♁ | 40% | | | | | | | | | | |

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.

We first examine results aggregated over difficulty conditions. The posterior predictives of all four RL-EAMs are shown in Figure 3, with the top row showing accuracies, and the middle and bottom rows correct and error RT distributions (parameter estimates for all models can be found in the supplementary materials). The RL-DDM generally explains the learning-related increase in accuracy well, and if only the central tendency were relevant it might be considered to provide an adequate account of RT, although correct median RT is systematically under-estimated. However, RT variability and skew are severely over-estimated. The RL-RD largely overcomes the RT distribution misfit, but it overestimates RTs in the first trial bins, and while capturing an increase in accuracy over trials, it is systematically underestimated. The RL-ARD

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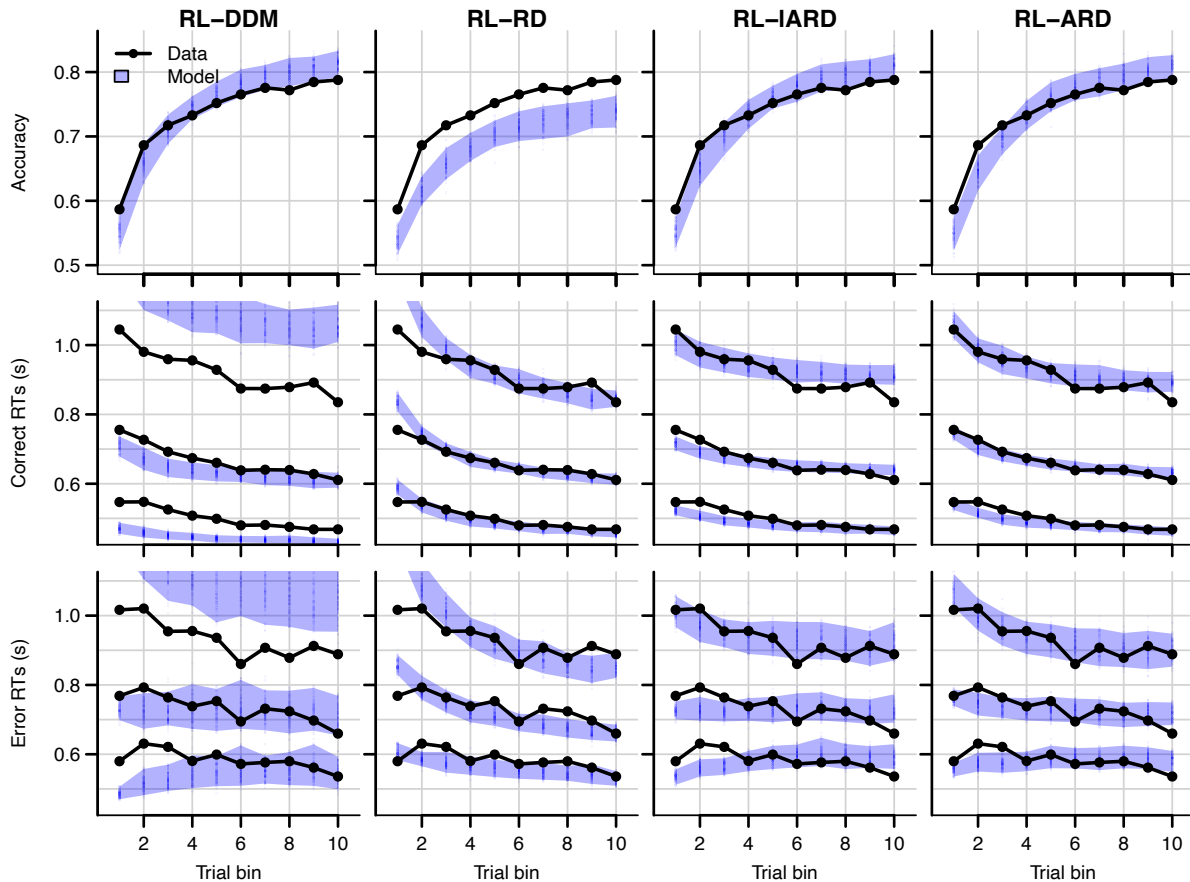179    except that the limited version slightly underestimates the decrease in RT with learning.
180



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-lARD, and RL-ARD (right column). Top row depicts accuracy over trial bins. Middle and bottom row show $10^{th}$, $50^{th,}$ and $90^{th}$ 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.
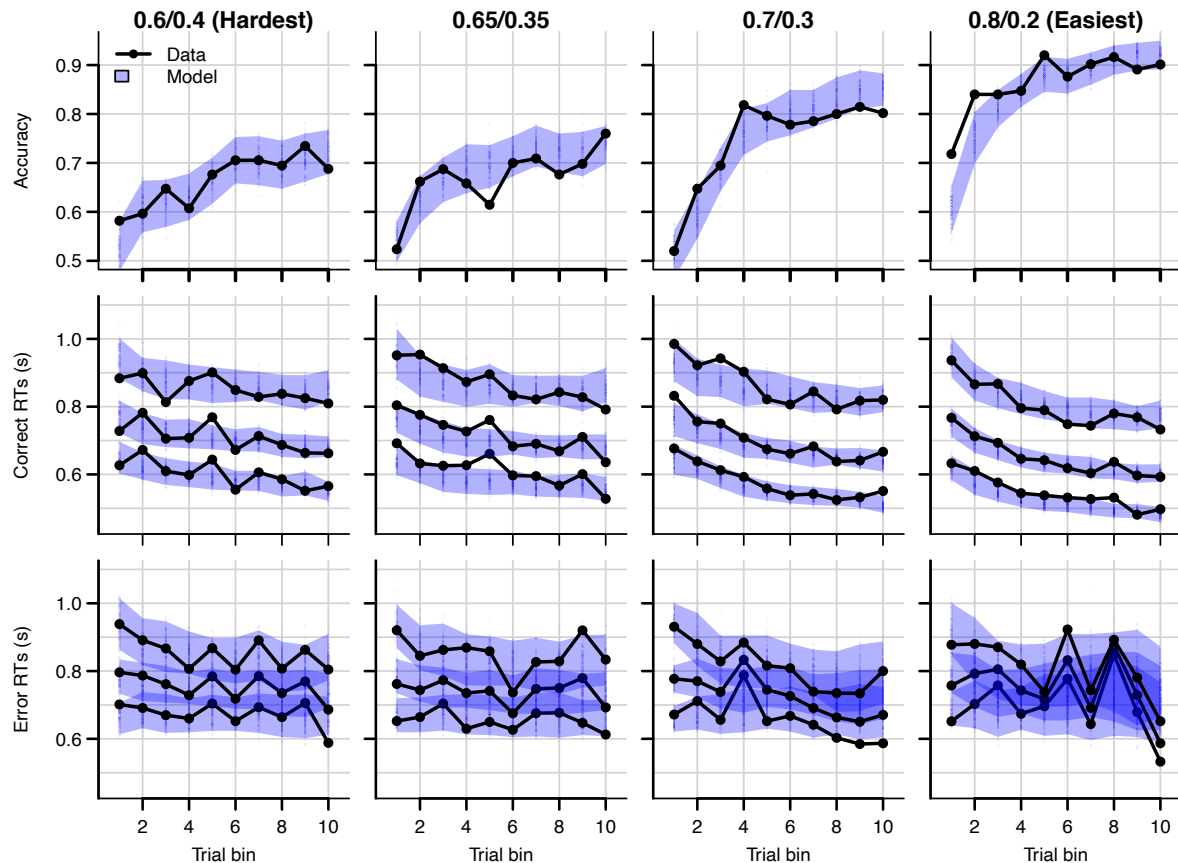
181
182
183

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 $10^{th}$, $50^{th}$, and $90^{th}$ 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 supplementary
186    materials for equivalent RL-DDM fits, which again fail to capture RT distributions). The RL-
187    ARD model displays the same excellent fit as to data aggregated over difficulty, except that it
188    underestimates accuracy in early trials in the easiest condition (Figure 4, bottom right panel).
189    Further inspections of the data revealed that 17 participants (31%) reached perfect accuracy in
190    the first bin in this condition. Likely, they guessed correctly on the first occurrence of the
191    easiest choice pair, repeated their choice, and received too little negative feedback in the next
192    repetitions to change their choice strategy. Supplementary materials show that, with these 17
193    participants removed, the overestimation is largely mitigated. SARSA assumes learning from
194    feedback, and so cannot explain such high early accuracies. Working memory processes could
195    have aided performance in the easiest condition, since the total number of stimuli pairs was
196    limited and feedback was quite reliable, making it relatively easy to remember correct-choice
197    options (Collins and Frank, 2012a, 2018; 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-lARD 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
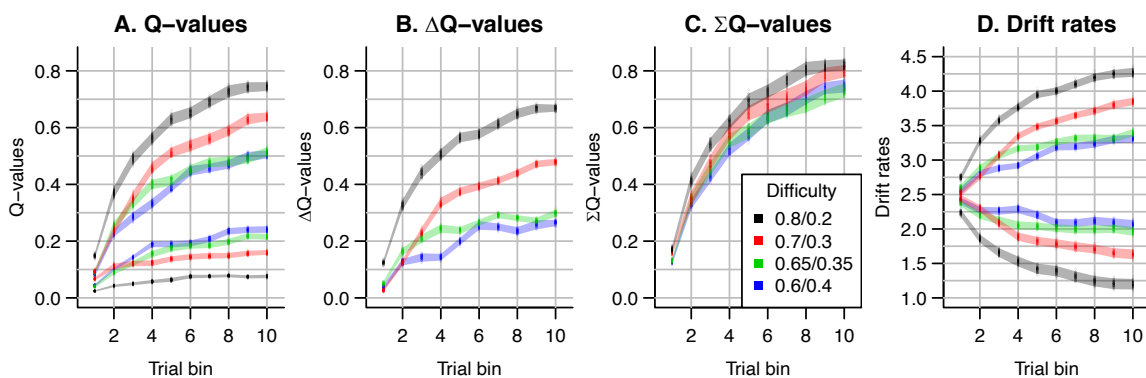218 effect on accuracy.



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

**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  Miletić, 2016; Miletić 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  To illustrate the importance of simultaneously analyzing RTs and choice behavior, we first
254  test whether a soft-max model (which ignores RTs) is able to capture the behavioral changes
255  in choice probability due to the manipulation. We fit two soft-max models to the data: One
256  with a single inverse temperature parameter, and one with an inverse temperature parameter
257  per SAT condition. The soft-max model with separate parameters per condition was
258  outperformed by a model with a single parameter ($\Delta BPIC = 11$), indicating that a researcher
259  using soft-max would have concluded that there was no difference in choice behavior between
260  conditions. Clearly, the difference in accuracy (and RTs) did indicate there were differences in
261  behavior (see supplementary materials for formal statistical tests), showing that soft-max fails
262  to capture a strong and well-known phenomenon of decision-making.

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

270  Formal model comparison (see Supplementary Table S1 for all BPIC values) indicated that
271  the RL-ARD model combining response caution and urgency effects provides the best
272  explanation of the data, in line with earlier research in non-learning contexts (Katsimpokis et
273  al., 2020; Miletić and Van Maanen, 2019; Rae et al., 2014; Thura and Cisek, 2016). The

274    advantage for the RL-ARD was substantial; the best RL-DDM (with only a threshold effect)
275    performed worse than the worst RL-ARD model. The data and posterior predictive
276    distributions of the best RL-DDM model and the winning RL-ARD model are shown in Figure
277    6. As in experiment 1, the RL-DDM failed to capture the shape of RT distributions, although
278    it fit the SAT effect on accuracy and median RTs. The RL-ARD model provides a much better
279    account of the RT distributions, including the differences between SAT conditions. In
280    supplementary materials we show that adding non-decision time variability to the RL-DDM
281    mitigates some of the misfit of the RT distributions, although it still consistently under-
282    predicted the $10^{th}$ percentile in the accuracy condition. Further, this model was still
283    substantially outperformed by the RL-ARD in formal model selection ($\Delta$BPIC = 209), and non-
284    decision time variability was estimated as much greater than what is found in non-learning
285    context, raising the question of its psychological plausibility.

286        Both RL-DDM and RL-ARD models tended to underestimate RTs and choice accuracy in
287    the early trial bins in the accuracy emphasis condition. As in experiment 1, working memory
288    may have contributed to the accurate but slow responses in the first trial bin for the accuracy
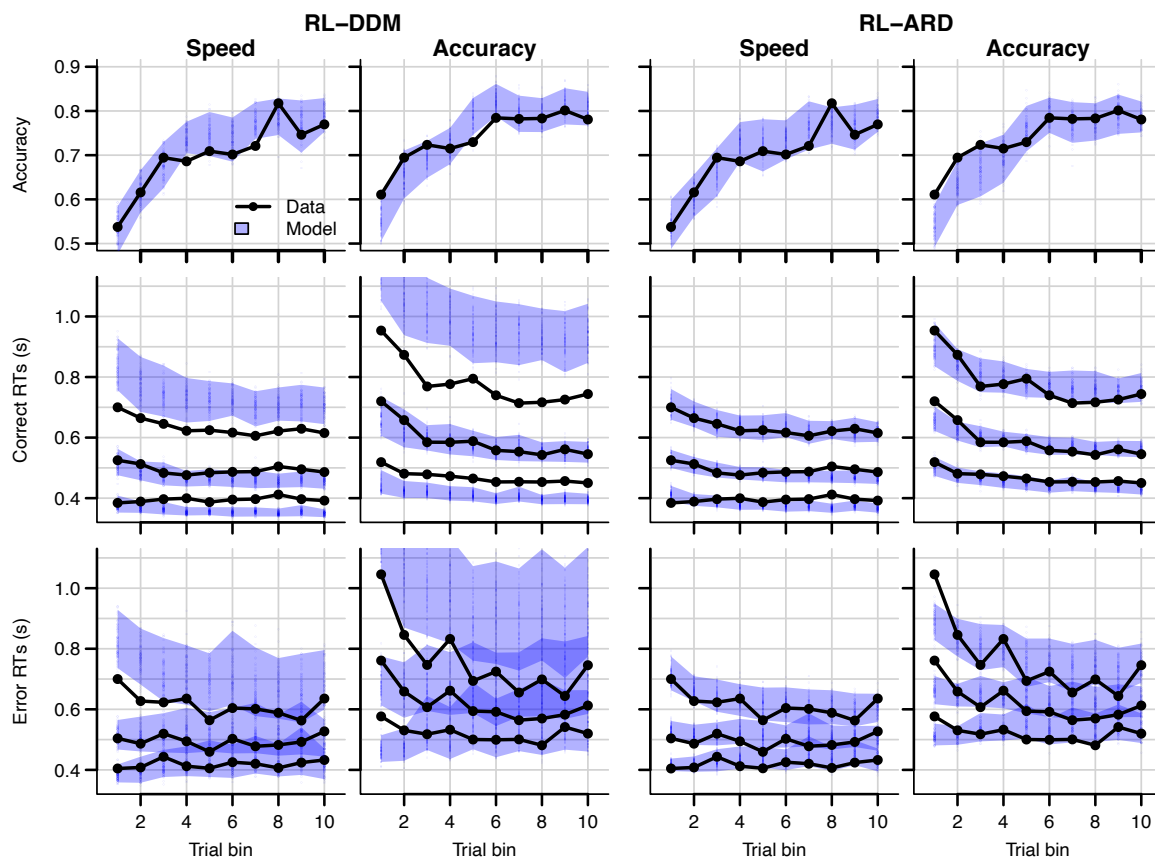289    condition (Collins and Frank, 2018, 2012b; McDougle and Collins, 2020).
290



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 $10^{th}$, $50^{th}$ and $90^{th}$ 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.

291

**Reversal learning**

292     Finally, we tested whether the RL-ARD can capture changes in accuracy and RTs caused by a
293     perturbation in the learning process due to reversals in reward contingencies. In the reversal
294     learning paradigm (Behrens et al., 2007; Costa et al., 2015; Izquierdo et al., 2017) participants
295     first learn a contingency between choice options and probabilistic rewards (the acquisition
296     phase) that is then suddenly reversed without any warning (the reversal phase). If the link
297     between Q-values and decision mechanisms as proposed by the RL-ARD underlies decisions,
298     the model should be able to account for the behavioral consequences (RT distributions and
299     decisions) of Q-value changes induced by the reversal.
300
301     Our reversal learning task had the same general structure as experiment 1 (Figure 1), except
302     for the presence of reversals. 47 participants completed four blocks of 128 trials each. Within
303     each block, two pairs of stimuli were randomly interleaved. Between trials 61 and 68
304     (uniformly sampled) in each block, the reward probability switched between stimuli, such that
305     stimuli that were correct during acquisition were incorrect after reversal (and vice versa).
306     Participants were not informed of the reversals prior to the experiment, but many reported
307     noticing them.
308     Data and the posterior predictive distributions of the RL-DDM and the RL-ARD models are
309     shown in Figure 7. Both models captured the change in choice proportions after the reversal
310     reasonably well, although they underestimate the speed of change. In supplementary materials
311     we show that the same is true for a standard soft-max model, suggesting that the learning rule
312     is the cause of this problem. Recent evidence indicates that, instead of only estimating expected
313     values of both choice options by error-driven learning, participants may additionally learn the
314     task structure, estimate the probability of a reversal occurring and adjust choice behavior
315     accordingly. Such a model-based learning strategy could increase the speed with which choice
316     behavior changes after a reversal (Costa et al., 2015; Izquierdo et al., 2017; Jang et al., 2015),
317     but as yet a learning rule that implements this strategy has not been developed.
318     The change in RT around the reversal was less marked than the change in choice probability.
319     Once again, the RL-DDM overestimates variability and skew. Both models fit the effects of
320     learning and reversal similarly, but the fastest responses for the RL-DDM decrease much too
321     quickly during initial learning and the reduction in speed for the slowest responses due to the
322     reversal is strongly overestimated. The RL-ARD provides a much better account of the shape
323     of the RT distributions, and furthermore captures the increase in entire RT *distributions*
324     (instead of only the median) after the reversal point. Formal model comparison also very
325     strongly favors the RL-ARD over the RL-DDM ($\Delta BPIC = 4051$). Supplementary materials
326     provide model comparisons to RL-DDMs with between-trial variability parameters, which lead
327     to the same conclusion.
328     A notable aspect of the data is that choice behavior stabilizes approximately 20 trials after
329     the reversal, whereas RTs remain high compared to just prior to the reversal point for up to ~40
330     trials. The RL-ARD explains this behavior through relatively high Q-values for the choice
331     option that was correct during the acquisition (but not reversal) phase (i.e., choice A). Figure
332     8 depicts the evolution of Q-values, Q-value differences and sums, and drift rates in the RL-
333     ARD model. The Q-values for both choice options increase until the reversal (Figure 8A), with
334     a much faster increase for $Q_A$. At the reversal $Q_A$ decreases and $Q_B$ increases, but as $Q_A$

13

335 decreases faster than $Q_B$ increases there is a temporary decrease in Q-value sums (Figure 8C).
336 After approximately 10 trials post-reversal, $Q_B$ is higher than for $Q_A$, which flips the sign of
337 the Q-value differences (Figure 8B). However, $Q_A$ *after* the reversal remains higher than the
338 $Q_B$ *before* the reversal, which causes the (absolute) Q-value differences to be lower after the
339 reversal than before. As a consequence, the drift rates for B after the reversal remain lower than
340 the drift rates for A before the reversal, which increases RT. Clearly, it is important to take
341 account of the sum of inputs to accumulators as well as the difference between them in order
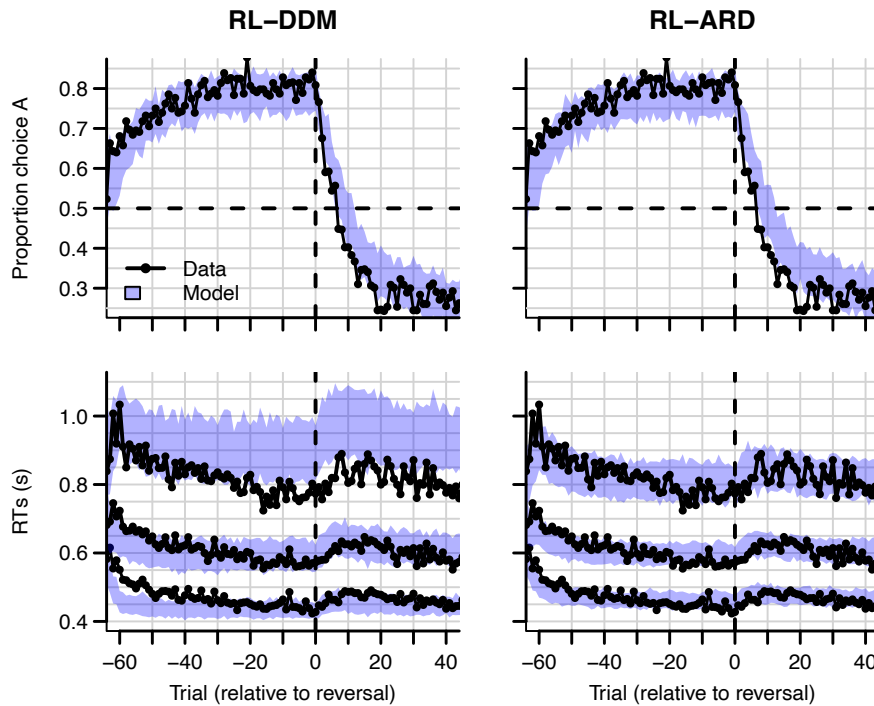342 to provide an accurate account of the effects of learning.

343



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.
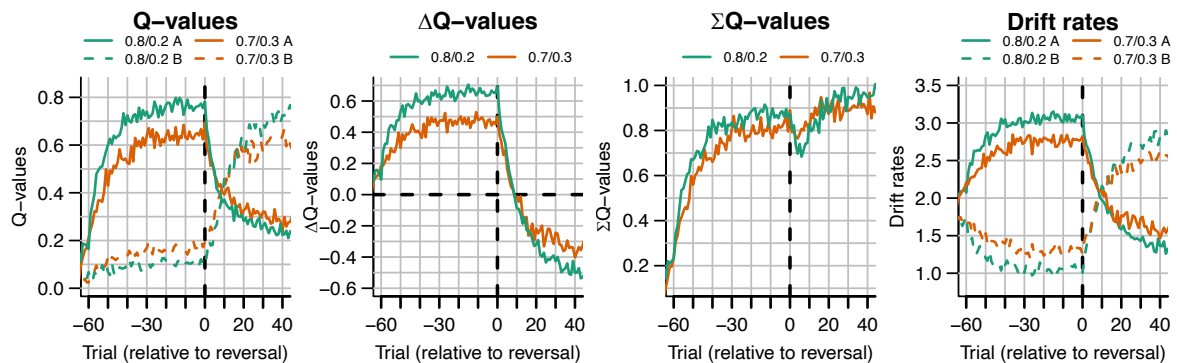
344



14

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.

345

346

**Discussion**

347 We compared combinations of different evidence-accumulations models with a simple SARSA
348 (Rummery and Niranjan, 1994) reinforcement learning rule (RL-EAMs). The comparison
349 tested the ability of the RL-EAMs to provide a comprehensive account of behavior in learning
350 contexts, not only in terms of the choices made but also the full distribution of the times to
351 make them (RT). We examined a standard instrumental learning paradigm (Frank, 2004) that
352 manipulated the difference in rewards between binary options (i.e., decision difficulty). We
353 also examined two elaborations of that paradigm testing key phenomena from the decision-
354 making and learning literatures, speed-accuracy trade-offs (SAT), and reward reversals,
355 respectively. Our benchmark was the dual threshold Diffusion Decision Model (Ratcliff, 1978)
356 (DDM), which has been used in almost all previous RL-EAM research, but has not been
357 compared to other RL-EAMs, and has not been thoroughly evaluated on its ability to account
358 for RT distributions in learning tasks. Our comparison used several different racing diffusion
359 (RD) models, where decisions depend on the winner of a race between single barrier diffusion
360 processes.
361
362     The RL-DDM provided a markedly inferior account to the other models, consistently
363 overestimating RT variability and skew. As these aspects of behavior are considered critical in
364 evaluating models in decision-making literature (Forstmann et al., 2016; Ratcliff and McKoon,
365 2008; Voss et al., 2013), our results question whether the RL-DDM provides an adequate
366 model of instrumental learning. Furthermore, the DDM carries with it two important theoretical
367 limitations. First, it can only address binary choice. This is unfortunate given that perhaps the
368 most widely used clinical application of reinforcement learning, the Iowa gambling task
369 (Bechara et al., 1994), requires choices among four options. Second, the input to the DDM
370 combines the evidence for each choice (i.e., "Q" values determined by the learning rule) into a
371 single difference, and so requires extra mechanisms to account for known effects of overall
372 reward magnitude (Fontanesi et al., 2019a). Although there are potential ways that the RL-
373 DDM might be modified to account for magnitude effects, such as increasing between-trial
374 drift rate variability in proportion to the mean rate (Ratcliff et al., 2018), its inability to extend
375 beyond binary choice remains an enduring impediment.
376     The best alternative model that we tested, the RL-ARD (advantage racing diffusion), which
377 is based on the recently proposed advantage accumulation framework (Van Ravenzwaaij et al.,
378 2020), remedied all of these problems. The input to each accumulator is the weighted sum of
379 three components: stimulus independent "urgency", the difference between evidence for the
380 choice corresponding to the accumulator and the alternative (the advantage), and the sum of
381 the evidence over accumulators. The urgency component had large effect in all fits and played

a key role in explaining the effect of speed-accuracy trade-offs. The advantage component, which is similar to the input to the DDM, was strongly supported over a model in which each accumulator only receives evidence favoring its own choice. The sum component provided a simple and theoretically transparent way to deal with reward magnitude effects in instrumental learning. Despite having the weakest effect among the three components, the sum was clearly necessary to provide an accurate fit to our data, even though we did not manipulate reward magnitude. It also played an important role in explaining the effect of reward reversals.

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

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

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

$$dx_1 = [V_0 + w_e Q_1 - w_i Q_2]dt + sW$$
$$dx_2 = [V_0 + w_e Q_2 - w_i Q_1]dt + sW$$

(5)

Where $w_e$ equals the sum of $w_D$ and $w_S$ in the parametrization of Equation (4), and $w_i$ equals the difference between $w_D$ and $w_S$. This re-parametrization shows that each drift rate is determined by an excitatory influence $w_e$ of the Q-value associated with the accumulator, and

424 an inhibitory influence $w_i$ of the Q-value associated with the other accumulator. Turner (2019)
425 proposed that inhibition plays an important role in learning tasks. Although the locus of
426 inhibition is different in the two models, there are clear parallels that bear further investigation.

427     A limitation of the current work is that we collapsed across blocks in analyzing the data of
428 experiments 2 and 3. However, in more detailed explorations (see supplementary materials for
429 details) there were indications of second-order changes across blocks. In experiment 2,
430 participants were faster in the first trial bin of the second and third block compared to the first
431 block, suggesting additional practice or adaptation effects at the beginning of the experiment.
432 In experiment 3, participants slowed down, and learned the reversal faster, after the first block.
433 This suggests they learned about the presence of reversals in the first block and applied a
434 different strategy in the later blocks. Although it is known that participants increase their
435 learning rates in volatile environments (Behrens et al., 2007), this by itself does not explain a
436 decrease in response speed. Potentially, if participants understood the task structure after the
437 first block, model-based strategies, such as estimating the probability of a reversal having
438 occurred, also slowed down responses.

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

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

463

464 **Methods**
465 **Experiment 1**
466 *Participants*

467    61 participants (mean age 21y [SD 2.33], 47 women, 56 right handed) were recruited from the
468    subject pool of the department of Psychology, University of Amsterdam, and participated for
469    course credits. All participants had normal or corrected-to-normal vision and gave written
470    informed consent prior to the experiment onset. The study was approved by the local ethics
471    committee.

472

473    *2.1.2 Task*
474    The task was an instrumental probabilistic learning task (Frank, 2004). On each trial, the
475    subject was presented with two abstract symbols (a 'stimulus pair') representing two choice
476    options (see Figure 2A for an example trial). Each choice option had a fixed probability of
477    being rewarded with points when chosen, with one choice option always having a higher
478    probability of being rewarded than the other. The task is to discover, by trial and error, which
479    choice options are most likely to lead to rewards, and thereby to collect as many points as
480    possible.
481        After a short practice block to get familiar with the task, participants completed one block
482    of 208 trials. Four different pairs of abstract symbols were included, each presented 52 times.
483    Stimulus pairs differed in their associated reward probabilities: 0.8/0.2, 0.7/0.3, 0.65/0.35, and
484    0.6/0.4. The size of the reward, if obtained, was always the same: '+100' (or '+0' otherwise).
485    Reward probabilities were chosen such that they differed only in the between-choice difference
486    in reward probability, leading to varying choice difficulties while keeping the mean reward
487    magnitude fixed.
488        Participants were instructed to earn as many points as possible, and to always respond before
489    the deadline of 2 seconds. Feedback consisted of two parts: an 'outcome' and a 'reward'. The
490    outcome corresponded to the probabilistic outcome of the choice, whereas the reward
491    corresponded to the actual number of earned points. When participants responded before the
492    deadline, the reward was equal to the outcome. If they were too late, the outcome was shown
493    to allow participants to learn from their choice, but the reward they received was set to 0 to
494    encourage responding in time. Participants received a bonus depending on the number of points
495    earned (maximum +0.5 course credits, mean received +0.24). The task was coded in PsychoPy
496    (Peirce et al., 2019). After this block, participants performed two more blocks of the same task
497    with different manipulations, which are not of current interest.

498

499    *Exclusion*
500    Six participants were excluded from analysis: One reported, after the experiment, not to have
501    understood the task, one reported a technical issue, and four did not reach an above-chance
502    accuracy level as determined by a binomial test (accuracy cut-off 0.55, corresponding to
503    $p < 0.05$). The final sample thus consisted of 55 subjects (14 men, mean age 21 years old [SD
504    2.39], 51 right-handed).

505

506    *Cognitive modelling*
507    The main analysis consists of fitting four RL-EAMs to the data and comparing the quality of
508    the fits penalized by model complexity. We compared four different decision models: the DDM
509    (Ratcliff, 1978), a racing diffusion (Boucher et al., 2007; Logan et al., 2014; Purcell et al.,
510    2010; Turner, 2019) model, and two Advantage Racing Diffusion (ARD; Van Ravenzwaaij et

511 al., 2020) models (see Figure 1 for an overview). Whereas the former is a two-sided diffusion
512 process, the latter three models employ a race architecture.
513     For all models we used the State-Action-Reward-State-Action (SARSA; Rummery and
514 Niranjan, 1994) update rule as a learning model:
515

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

516

517 where $Q_{i,t}$ is the value representation of choice option $i$ on trial $t$, $\alpha$ the learning rate, and $r_t$
518 the reward on trial $t$. The difference between the actual reward and the value representation of
519 the chosen stimulus, $r_t - Q_{i,t}$, is known as the reward prediction error. The learning rate
520 controls the speed at which Q-values change in response to the reward prediction error, with
521 larger learning rates leading to stronger fluctuations. In this model, only the Q-value of the
522 chosen option is updated.
523

524 *RL-EAM 1: RL-DDM*
525 In the first RL-EAM, we use the DDM (Ratcliff, 1978) as a choice model (Figure 1, left
526 column). The DDM assumes that evidence accumulation is governed by:
527

$$dx = vdt + sW$$

528

529 $v$ is the mean speed of evidence accumulation (the *drift rate*), and $s$ is the standard deviation
530 of the within-trial accumulation white noise (W). The RL-DDM assumes that the drift rate
531 depends linearly on the difference of value representations:
532

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

533

534 $w$ is a weighting variable, and $Q_{t,1}$ and $Q_{t,2}$ are the $Q$-values for both choice options per trial,
535 which change each trial according to Equation 4. Hence,
536

$$dx = w(Q_1 - Q_2)\, dt + sW \tag{1}$$

537

538 The starting point of evidence accumulation, $z$, lies between decision boundaries $a$ and $-a$.
539 Here, as in earlier RL-DDM work (Fontanesi et al., 2019a, 2019b; Pedersen et al., 2017), we
540 assume an unbiased start of the decision process (i.e., $z = 0$). Evidence accumulation finishes
541 when threshold $a$ or $-a$ is reached, and the decision for the choice corresponding to $Q_1$ or $Q_2$,
542 respectively, is made. The response time is the time required for the evidence-accumulation
543 process to reach the bound, plus an intercept called the non-decision time ($t0$). The non-
544 decision time is the sum of the time required for perceptual encoding and the time required for
545 the execution of the motor response. Parameter $s$ was fixed to 1 to satisfy scaling constraints
546 (Donkin et al., 2009; van Maanen and Miletić, 2020). In total, this specification of the RL-
547 DDM has 4 free parameters ($\alpha, w, a, t0$). In the supplementary materials, we fit an RL-DDM
548 specification with between-trial variabilities in start point, drift rate, and non-decision time.

549
550   *RL-EAM 2: RL-RD*
551   The RL-RD (Figure 1, middle panel) assumes that two evidence accumulators independently
552   accrue evidence for one choice option each, both racing towards a common threshold $a$
553   (assuming no response bias). The first accumulator to hit the bound wins, and the
554   corresponding decision is made. For each choice option $i$, the dynamics of accumulation are
555   governed by:
556

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

557
558   $V_0$ is a parameter specifying the drift rate in the absence of any evidence, $w$ a weighting
559   parameter, and $s$ the standard deviation of within-trial noise. As such, the mean speed of
560   accumulation (the drift rate $v_i$) is the sum of two independent factors: an evidence-independent
561   baseline speed $V_0$, and an evidence-dependent weighted Q-value, $wQ_i$. Since $V_0$ is assumed to
562   be identical across accumulators, and governs the speed of accumulation unrelated to the
563   amount of evidence, we interpret this parameter as an additive urgency signal (Miletić and Van
564   Maanen, 2019), with conceptually similar behavioral effects as collapsing bounds (Hawkins et
565   al., 2015). Similar to the DDM, a non-decision time parameter accounts for the time for
566   perceptual encoding and the motor response time. Parameter $s$ was fixed to 1 to satisfy scaling
567   constraints (Donkin et al., 2009; van Maanen and Miletić, 2020). In total, the RL-RD has 5 free
568   parameters $(\alpha, w, a, v0, t0)$.
569      Each accumulator's first passage times are Wald (also known as inverted Gaussian)
570   distributed (Anders et al., 2016). In an independent race model, each accumulator's first
571   passage time distribution is normalized to the probability of the response with which it is
572   associated (Brown and Heathcote, 2008; Turner, 2019).
573
574   *RL-EAM 3 & 4: RL-ARD*
575   Thirdly, we fit two racing diffusion models based on an advantage race architecture (Van
576   Ravenzwaaij et al., 2020). An advantage race model using an LBA has been shown to provide
577   a natural account for multi-alternative choice phenomena such as Hick's law, as well as
578   stimulus magnitude effects in perceptual decision-making. Like in the RL-RD, accumulators
579   race towards a common bound, but the speed of evidence accumulation $v_i$ depends on multiple
580   factors: first, as in the RL-RD, the evidence-independent speed of accumulation $V_0$; second,
581   the *advantage* of the evidence for one choice option over the other (c.f. the DDM, where the
582   difference between evidence for both choice options is accumulated); and third, the *sum* of the
583   total available evidence. Combined, for two accumulators in the RL-EAM framework, this
584   leads to:
585

$$dx_1 = [V_0 + w_d(Q_1 - Q_2) + w_s(Q_1 + Q_2)]dt + sW$$
$$dx_2 = [V_0 + w_d(Q_2 - Q_1) + w_s(Q_1 + Q_2)]dt + sW \qquad (3)$$

586
587   In the original work proposing the advantage accumulation framework (Van Ravenzwaaij et
588   al., 2020), it was shown that the $w_d$ parameter had a much stronger influence on evidence-

589 accumulation rates than the $w_s$ parameter. Therefore, we first fixed the $w_s$ parameter to 0, to
590 test whether the accumulation of *differences* is sufficient to capture all trends in the data. We
591 term this model the RL-lARD (l = limited), which we compare to the RL-ARD in which we fit
592 $w_s$ as a free parameter.

593    As previously, parameter *s* was fixed to 1 to satisfy scaling constraints (Donkin et al., 2009;
594 van Maanen and Miletić, 2020). The RL-ARD also has a threshold, non-decision time, and
595 learning rate parameter, totaling five $(\alpha, w_d, a, V_0, t0)$ and 6 free parameters
596 $(\alpha, w_d, w_s, a, V_0, t0)$ for the RL-lARD and RL-ARD, respectively.

597

598 *Bayesian hierarchical parameter estimation, posterior predictive distributions, model*
599 *comparisons*
600 We estimated group-level and subject-level posterior distributions of each model's parameter
601 using a combination of differential evolution (DE) and Markov-chain Monte Carlo sampling
602 (MCMC) with Metropolis-Hastings (Ter Braak, 2006; Turner et al., 2013). Sampling settings
603 were default as implemented in the Dynamic Models of Choice *R* software (Heathcote et al.,
604 2019): The number of chains, D, was three times the number of free parameters. Cross-over
605 probability was set to $2.38/\sqrt{D}$ at the subject-level and $U[0,1]$ at the group level. Migration
606 probability was set to 0.05 during burn-in only. Convergence was assessed using visual
607 inspection of the chain traces and Gelman-Rubin diagnostic (Brooks and Gelman, 1998;
608 Gelman and Rubin, 1992) (individual and multivariate potential scale factors < 1.03 in all
609 cases).
610    Hierarchical models were fit assuming independent normal population ("hyper")
611 distributions for each parameter. For all models, we estimated the learning rate on a probit scale
612 (mapping [0, 1] onto the real domain), with a normal prior $\alpha \sim \Phi(\mathcal{N}(-1.6, 5))$ (Spektor and
613 Kellen, 2018). Prior distributions for all estimated hyper-mean decision-related parameters
614 were vague. RL-EAMs, the threshold parameter $a \sim \mathcal{N}(3, 5)$ truncated at 0, and
615 $t0 \sim \mathcal{N}(0.3, 0.5)$ truncated at 0.025 s and 1 s (all estimation was carried out on the seconds
616 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
617 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
618 as prior. Plots of superimposed prior and posterior hyper-distributions confirmed that these
619 prior setting were not influential.
620    In initial explorations, we also freely estimated the Q-values at trial 0. However, in the RL-
621 EAMs, the posterior distributions for these Q-values consistently converged on 0, which was
622 therefore subsequently used as a fixed value for all results reported here. For the the soft-max
623 fits, they were set to 0.5 as often used in reinforcement learning models of two-choice tasks
624 (Apps et al., 2015; Collins and Frank, 2018, 2012b; Fontanesi et al., 2019a; McDougle and
625 Collins, 2020; Pedersen and Frank, 2020). Including the initial Q-values as a free parameter in
626 the soft-max models of experiment 2 led to the same conclusions.
627    To visualize the quality of model fit, we took 100 random samples from the estimated
628 parameter posteriors and simulated the experimental design with these parameters. For each
629 behavioral measure (e.g., RT quantiles, accuracy), credible intervals were estimated by taking
630 the range between the 2.5% and 97.5% quantiles of the averages over participants.

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

638

639    **Experiment 2**
640    *Participants*
641    23 participants (mean age 19 years old [SD 1.06 years], 7 men, 23 right-handed) were recruited
642    from the subject pool of the Department of Psychology of the University of Amsterdam and
643    participated for course credits. Participants did not participate in experiment 1 or 3. All
644    participants had normal or corrected-to-normal vision and gave written informed consent prior
645    to the experiment onset. The study was approved by the local ethics committee.

646

647    *Task*
648    Participants performed the same task as in experiment 1, with the addition of an SAT
649    manipulation (Figure 2C). The SAT manipulation included both an instructional cue and a
650    response deadline. Prior to each trial, a cue instructed participants to emphasize either decision
651    speed ('SPD') or decision accuracy ('ACC') in the upcoming trial, and in speed trials,
652    participants did not earn points if they were too late (> 700 ms). As in experiment 1, after each
653    choice participants received feedback consisting of two components: an outcome and a reward.
654    The outcome refers to the outcome of the probabilistic gamble, whereas the reward refers to
655    the number of points participants actually received. If participants responded in time, the
656    reward was equal to the outcome. In speed trials, participants did not earn points if they
657    responded later than 700 ms after stimulus onset, even if the outcome was +100. On trials
658    where participants responded too late, they were additionally informed of the reward that was
659    associated with their choice, had they been in time. This way, even when participants are too
660    late, they still receive the feedback that can be used to learn from their choices.
661    The deadline manipulation was added because we hypothesized that instructional cues alone
662    would not be sufficient to persuade participants to change their behavior in the instrumental
663    learning task, since that task specifically requires them to accumulate points. If the received
664    number of points was independent of response times, the optimal strategy to collect most points
665    would be to ignore the cue and focus on accuracy only.
666    Participants performed 324 trials divided over 3 blocks. Within each block, three pairs of
667    stimuli were shown, with associated reward probabilities of 0.8/0.2, 0.7/0.3, and 0.6/0.4. Speed
668    and accuracy trials were randomly interleaved. Figure 2C depicts the sequence of events in
669    each trial. As this experiment also served as a pilot for an fMRI experiment, we added fixation
670    crosses between each phase of the trial, with jittered durations. A pre-stimulus fixation cross
671    lasted 0.5, 1, 1.5, or 2 s; fixation crosses between cue and stimulus, between stimulus and
672    highlight, and between highlight and feedback lasted 0, 0.5, 1, or 1.5 s; and an inter-trial
673    interval fixation cross lasted 0.5, 1, 1.5, 2, 2.5 seconds. Each trial took 7.5 seconds. The
674    experiment took approximately 45 minutes.

*Exclusion*

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

*Cognitive modelling*

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

$$P_{i,t} = \frac{\exp \beta Q_{i,t}}{\sum_{j}^{J} \exp \beta Q_{j,t}} \qquad (5)$$

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

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

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

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

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

**Experiment 3**

*Participants*

713 47 participants (mean age 21 years old [SD 2.81 years], 16 men, 40 right-handed) were
714 recruited from the subject pool of the Department of Psychology of the University of
715 Amsterdam and participated for course credits. Participants did not participate in experiment 1
716 or 2. All participants had normal or corrected-to-normal vision and gave written informed
717 consent prior to the experiment onset. The study was approved by the local ethics committee.
718

719 *Task*
720 The reversal learning task had the same general task structure as experiment 1. Participants
721 completed four blocks of 128 trials each, totaling 512 trials. Within each block, two pairs of
722 stimuli were randomly interleaved, with associated reward probabilities of 0.8/0.2 and 0.7/0.3.
723 Between trials 61 and 68 (uniformly sampled) of each block, the reward probability switched
724 between stimuli, such that the stimulus with a pre-reversal reward probability of 0.8/0.7 had a
725 post-reversal reward probability of 0.2/0.3 (and vice versa). Participants were not informed of
726 the reversals prior to the experiment, but many reported noticing them.
727 In addition to the reversal learning task, the experimental session also contained a working
728 memory task that is not of current interest. 30 participants performed the reversal learning task
729 before the working memory task, and 17 participants afterwards. The entire experiment took
730 approximately one hour.
731

732 *Cognitive modelling*
733 The RL-DDM and RL-ARD were fit to the data using the same methods as in experiment 1.
734

735

736 **Data availability statement**
737 All data are available on OSF (https://osf.io/ygrve/).
738
739 **Code availability statement**
740 All analysis code is available on OSF (https://osf.io/ygrve/).
741

746

747

748

749 **References**
750 Anders R, Alario F, Van Maanen L. 2016. The Shifted Wald Distribution for Response Time Data Analysis.
751 *Psychol Methods* **21**:309–327.
752 Ando T. 2007. Bayesian predictive information criterion for the evaluation of hierarchical Bayesian and
753 empirical Bayes models. *Biometrika* **94**:443–458. doi:10.1093/biomet/asm017
754 Apps MAJ, Lesage E, Ramnani N. 2015. Vicarious reinforcement learning signáis when instructing others. *J
755 Neurosci* **35**:2904–2913. doi:10.1523/JNEUROSCI.3669-14.2015
756 Arnold NR, Bröder A, Bayen UJ. 2015. Empirical validation of the diffusion model for recognition memory and
757 a comparison of parameter-estimation methods. *Psychol Res* **79**:882–898. doi:10.1007/s00426-014-0608-y
758 Barto AG, Sutton RS, Brouwer PS. 1981. Associative search network: A reinforcement learning associative
759 memory. *Biol Cybern* **40**:201–211. doi:10.1007/BF00453370

Bechara a, Damasio a R, Damasio H, Anderson SW. 1994. Insensitivity to future consequences following damage to human prefrontal cortex. *Cognition* **50**:7–15. doi:10.1016/0010-0277(94)90018-3

Behrens TEJ, Woolrich MW, Walton ME, Rushworth MFS. 2007. Learning the value of information in an uncertain world. *Nat Neurosci* **10**:1214–1221. doi:10.1038/nn1954

Boag RJ, Strickland L, Heathcote A, Neal A, Loft S. 2019a. Cognitive Control and Capacity for Prospective Memory in Complex Dynamic Environments. *J Exp Psychol Gen* **148**:2181–2206. doi:10.1037/xge0000599

Boag RJ, Strickland L, Loft S, Heathcote A. 2019b. Strategic attention and decision control support prospective memory in a complex dual-task environment. *Cognition* **191**:103974. doi:10.1016/j.cognition.2019.05.011

Bogacz R, Larsen T. 2011. Integration of reinforcement learning and optimal decision-making theories of the basal ganglia. *Neural Comput* **23**:817–851. doi:10.1162/NECO_a_00103

Bogacz R, McClure SM, Li J, Cohen JD, Montague PR. 2007. Short-term memory traces for action bias in human reinforcement learning. *Brain Res* **1153**:111–121. doi:10.1016/j.brainres.2007.03.057

Bogacz R, Wagenmakers E-J, Forstmann BU, Nieuwenhuis S. 2010. The neural basis of the speed-accuracy tradeoff. *Trends Neurosci* **33**:10–6. doi:10.1016/j.tins.2009.09.002

Boucher L, Palmeri TJ, Logan GD, Schall JD. 2007. Inhibitory Control in Mind and Brain : An Interactive Race Model of Countermanding Saccades **114**:376–397. doi:10.1037/0033-295X.114.2.376

Brooks SP, Gelman A. 1998. General Methods for Monitoring Convergence of Iterative Simulations. *J Comput Graph Stat* **7**:434–455. doi:10.1080/10618600.1998.10474787

Brown SD, Heathcote A. 2008. The simplest complete model of choice response time: Linear ballistic accumulation. *Cogn Psychol* **57**:153–178. doi:10.1016/j.cogpsych.2007.12.002

Christakou A, Gershman SJ, Niv Y, Simmons A, Brammer M, Rubia K. 2013. Neural and Psychological Maturation of Decision-making in Adolescence and Young Adulthood. *J Cogn Neurosci* **25**:1807–1823. doi:10.1162/jocn_a_00447

Cisek P, Puskas GA, El-Murr S. 2009. Decisions in changing conditions: the urgency-gating model. *J Neurosci* **29**:11560–71. doi:10.1523/JNEUROSCI.1844-09.2009

Collins AGE, Frank MJ. 2018. Within- and across-trial dynamics of human EEG reveal cooperative interplay between reinforcement learning and working memory. *Proc Natl Acad Sci* **115**:2502–2507. doi:10.1073/pnas.1720963115

Collins AGE, Frank MJ. 2012a. How much of reinforcement learning is working memory, not reinforcement learning? A behavioral, computational, and neurogenetic analysis. *Eur J Neurosci* **35**:1024–1035. doi:10.1111/j.1460-9568.2011.07980.x

Collins AGE, Frank MJ. 2012b. How much of reinforcement learning is working memory, not reinforcement learning? A behavioral, computational, and neurogenetic analysis. *Eur J Neurosci* **35**:1024–1035. doi:10.1111/j.1460-9568.2011.07980.x

Costa VD, Tran VL, Turchi J, Averbeck BB. 2015. Reversal learning and dopamine: A Bayesian perspective. *J Neurosci* **35**:2407–2416. doi:10.1523/JNEUROSCI.1989-14.2015

Daw ND, Dayan P. 2014. The algorithmic anatomy of model-based evaluation. *Philos Trans R Soc B Biol Sci* **369**. doi:10.1098/rstb.2013.0478

Daw ND, Kakade S, Dayan P. 2002. Opponent interactions between serotonin and dopamine. *Neural Networks* **15**:603–616. doi:10.1016/S0893-6080(02)00052-7

Daw ND, O'Doherty JP, Dayan P, Seymour B, Dolan RJ. 2006. Cortical substrates for exploratory decisions in humans. *Nature* **441**:876–879. doi:10.1038/nature04766

Dayan P, Daw ND. 2008. Decision theory, reinforcement learning, and the brain. *Cogn Affect Behav Neurosci* **8**:429–453. doi:10.3758/CABN.8.4.429

Donkin C, Brown SD. 2018. Response Times and Decision-Making, Stevens' Handbook of Experimental Psychology and Cognitive Neuroscience. doi:10.1002/9781119170174.epcn509

Donkin C, Brown SD, Heathcote A. 2011. Drawing conclusions from choice response time models: A tutorial using the linear ballistic accumulator. *J Math Psychol* **55**:140–151. doi:10.1016/j.jmp.2010.10.001

Donkin C, Brown SD, Heathcote A. 2009. The overconstraint of response time models: Rethinking the scaling problem. *Psychon Bull Rev* **16**:1129–1135. doi:10.3758/PBR.16.6.1129

Dutilh G, Rieskamp J. 2016. Comparing perceptual and preferential decision making. *Psychon Bull Rev* **23**:723–737. doi:10.3758/s13423-015-0941-1

Fontanesi L, Gluth S, Spektor MS, Rieskamp J. 2019a. A reinforcement learning diffusion decision model for value-based decisions. *Psychon Bull Rev*. doi:10.3758/s13423-018-1554-2

Fontanesi L, Palminteri S, Lebreton M. 2019b. Decomposing the effects of context valence and feedback information on speed and accuracy during reinforcement learning: a meta-analytical approach using diffusion decision modeling. *Cogn Affect Behav Neurosci* **19**:490–502. doi:10.3758/s13415-019-00723-1

Forstmann BU, Ratcliff R, Wagenmakers E-J. 2016. Sequential Sampling Models in Cognitive Neuroscience: Advantages, Applications, and Extensions. *Annu Rev Psychol* **67**:641–666. doi:10.1146/annurev-psych-

122414-033645

Frank MJ. 2004. By Carrot or by Stick: Cognitive Reinforcement Learning in Parkinsonism. *Science (80- )* **306**:1940–1943. doi:10.1126/science.1102941

Frank MJ, Doll BB, Oas-Terpstra J, Moreno F. 2009. Prefrontal and striatal dopaminergic genes predict individual differences in exploration and exploitation. *Nat Neurosci* **12**:1062–1068. doi:10.1038/nn.2342

Gelman A, Rubin DB. 1992. Inference from Iterative Simulation Using Multiple Sequences. *Stat Sci* **7**:457–472. doi:10.1214/ss/1177011136

Gershman SJ. 2015. Do learning rates adapt to the distribution of rewards? *Psychon Bull Rev* **22**:1320–1327. doi:10.3758/s13423-014-0790-3

Haughey HM, Hutchison KE, Curran T, Frank MJ, Moustafa AA. 2007. Genetic triple dissociation reveals multiple roles for dopamine in reinforcement learning. *Proc Natl Acad Sci* **104**:16311–16316. doi:10.1073/pnas.0706111104

Hawkins GE, Forstmann BU, Wagenmakers E-J, Ratcliff R, Brown SD. 2015. Revisiting the Evidence for Collapsing Boundaries and Urgency Signals in Perceptual Decision-Making. *J Neurosci* **35**:2476–2484. doi:10.1523/JNEUROSCI.2410-14.2015

Hawkins GE, Heathcote A. 2020. Racing Against The Clock: Evidence-Based Vs. Time-Based Decisions. *Psychol Rev*.

Heathcote A, Lin YS, Reynolds A, Strickland L, Gretton M, Matzke D. 2019. Dynamic models of choice. *Behav Res Methods* **51**:961–985. doi:10.3758/s13428-018-1067-y

Heathcote A, Love J. 2012. Linear deterministic accumulator models of simple choice. *Front Psychol* **3**:1–19. doi:10.3389/fpsyg.2012.00292

Ho T, Brown SD, Van Maanen L, Forstmann BU, Wagenmakers E-J, Serences JT. 2012. The Optimality of Sensory Processing during the Speed-Accuracy Tradeoff. *J Neurosci* **32**:7992–8003. doi:10.1523/JNEUROSCI.0340-12.2012

Izquierdo A, Brigman JL, Radke AK, Rudebeck PH, Holmes A. 2017. The neural basis of reversal learning: An updated perspective. *Neuroscience* **345**:12–26. doi:10.1016/j.neuroscience.2016.03.021

Jang AI, Costa VD, Rudebeck PH, Chudasama Y, Murray EA, Averbeck BB. 2015. The role of frontal cortical and medial-temporal lobe brain areas in learning a Bayesian prior belief on reversals. *J Neurosci* **35**:11751–11760. doi:10.1523/JNEUROSCI.1594-15.2015

Katsimpokis D, Hawkins GE, Van Maanen L. 2020. Not all Speed-Accuracy Trade-Off Manipulations Have the Same Psychological Effect. *Comput Brain Behav*. doi:10.1007/s42113-020-00074-y

Leite FP, Ratcliff R. 2010. Modeling reaction time and accuracy of multiple-alternative decisions. *Attention, Perception, Psychophys* **72**:246–273. doi:10.3758/APP.72.1.246

Logan GD, Van Zandt T, Verbruggen F, Wagenmakers EJ. 2014. On the ability to inhibit thought and action: General and special theories of an act of control. *Psychol Rev* **121**:66–95. doi:10.1037/a0035230

Luzardo A, Alonso E, Mondragón E. 2017. A Rescorla-Wagner drift-diffusion model of conditioning and timing. *PLOS Comput Biol* **13**:e1005796. doi:10.1371/journal.pcbi.1005796

McDougle SD, Collins AGE. 2020. Modeling the influence of working memory, reinforcement, and action uncertainty on reaction time and choice during instrumental learning. *Psychon Bull Rev*. doi:10.3758/s13423-020-01774-z

Miletić S. 2016. Neural Evidence for a Role of Urgency in the Speed-Accuracy Trade-off in Perceptual Decision-Making. *J Neurosci* **36**:5909–5910. doi:10.1523/JNEUROSCI.0894-16.2016

Miletić S, Boag RJ, Forstmann BU. 2020. Mutual benefits: Combining reinforcement learning with sequential sampling models. *Neuropsychologia* **136**. doi:10.1016/j.neuropsychologia.2019.107261

Miletić S, Van Maanen L. 2019. Caution in decision-making under time pressure is mediated by timing ability. *Cogn Psychol* **110**:16–29. doi:10.1016/j.cogpsych.2019.01.002

Millner AJ, Gershman SJ, Nock MK, den Ouden HEM. 2018. Pavlovian Control of Escape and Avoidance. *J Cogn Neurosci* **30**:1379–1390. doi:10.1162/jocn_a_01224

Murphy PR, Boonstra E, Nieuwenhuis S. 2016. Global gain modulation generates time-dependent urgency during perceptual choice in humans. *Nat Commun* **7**:1–14. doi:10.1038/ncomms13526

Niv Y, Edlund JA, Dayan P, O'Doherty JP. 2012. Neural Prediction Errors Reveal a Risk-Sensitive Reinforcement-Learning Process in the Human Brain. *J Neurosci* **32**:551–562. doi:10.1523/jneurosci.5498-10.2012

O'Doherty JP, Cockburn J, Pauli WM. 2017. Learning, Reward, and Decision Making. *Annu Rev Psychol* **68**:73–100. doi:10.1146/annurev-psych-010416-044216

Pachella RG, Pew RW. 1968. Speed-Accuracy Tradeoff in Reaction Time: Effect of Discrete Criterion Times. *J Exp Psychol* **76**:19–24. doi:10.1037/h0021275

Palminteri S, Khamassi M, Joffily M, Coricelli G. 2015. Contextual modulation of value signals in reward and punishment learning. *Nat Commun* **6**. doi:10.1038/ncomms9096

Palminteri S, Wyart V, Koechlin E. 2017. The Importance of Falsification in Computational Cognitive

Modeling. *Trends Cogn Sci* **21**:425–433. doi:10.1016/j.tics.2017.03.011

Pedersen ML, Frank MJ. 2020. Simultaneous Hierarchical Bayesian Parameter Estimation for Reinforcement Learning and Drift Diffusion Models: a Tutorial and Links to Neural Data. *Comput Brain Behav*. doi:10.1007/s42113-020-00084-w

Pedersen ML, Frank MJ, Biele G. 2017. The drift diffusion model as the choice rule in reinforcement learning. *Psychon Bull Rev* **24**:1234–1251. doi:10.3758/s13423-016-1199-y

Peirce J, Gray JR, Simpson S, MacAskill M, Höchenberger R, Sogo H, Kastman E, Lindeløv JK. 2019. PsychoPy2: Experiments in behavior made easy. *Behav Res Methods* **51**:195–203. doi:10.3758/s13428-018-01193-y

Purcell BA, Heitz RP, Cohen JY, Schall JD, Logan GD, Palmeri TJ. 2010. Neurally constrained modeling of perceptual decision making. *Psychol Rev* **117**:1113–1143. doi:10.1037/a0020311

Rae B, Heathcote A, Donkin C, Averell L, Brown SD. 2014. The hare and the tortoise: Emphasizing speed can change the evidence used to make decisions. *J Exp Psychol Learn Mem Cogn* 1–39. doi:10.1037/a0036801

Ratcliff R. 1978. A theory of memory retrieval. *Psychol Rev* **85**:59–108.

Ratcliff R, Hasegawa YT, Hasegawa RP, Childers R, Smith PL, Segraves MA. 2011. Inhibition in superior colliculus neurons in a brightness discrimination task? *Neural Comput* **23**:1790–1820. doi:10.1162/NECO_a_00135

Ratcliff R, Hasegawa YT, Hasegawa RP, Smith PL, Segraves MA. 2007. Dual Diffusion Model for Single-Cell Recording Data From the Superior Colliculus in a Brightness-Discrimination Task. *J Neurophysiol* **97**:1756–1774. doi:10.1152/jn.00393.2006

Ratcliff R, McKoon G. 2008. The diffusion decision model: theory and data for two-choice decision tasks. *Neural Comput* **20**:873–922. doi:10.1162/neco.2008.12-06-420

Ratcliff R, Rouder JN. 1998. Modeling Response Times for Two-Choice Decisions. *Psychol Sci* **9**:347–356.

Ratcliff R, Smith PL, Brown SD, McKoon G. 2016. Diffusion Decision Model: Current Issues and History. *Trends Cogn Sci* **20**:260–281. doi:10.1016/j.tics.2016.01.007

Ratcliff R, Voskuilen C, Teodorescu A. 2018. Modeling 2-alternative forced-choice tasks: Accounting for both magnitude and difference effects. *Cogn Psychol* **103**:1–22. doi:10.1016/j.cogpsych.2018.02.002

Rescorla RA, Wagner AR. 1972. A theory of Pavlovian conditioning: Variations in the effectiveness of reinforcement and nonreinforcement. *Class Cond II Curr Res Theory* **21**:64–99. doi:10.1101/gr.110528.110

Rummery GA, Niranjan M. 1994. On-Line Q-Learning Using Connectionist Systems.

Sewell DK, Jach HK, Boag RJ, Van Heer CA. 2019. Combining error-driven models of associative learning with evidence accumulation models of decision-making. *Psychon Bull Rev*. doi:10.3758/s13423-019-01570-4

Sewell DK, Stallman A. 2020. Modeling the Effect of Speed Emphasis in Probabilistic Category Learning. *Comput Brain Behav* **3**:129–152. doi:10.1007/s42113-019-00067-6

Shahar N, Hauser TU, Moutoussis M, Moran R, Keramati M, Consortium NSPN, Dolan RJ. 2019. Improving the reliability of model-based decision-making estimates in the two-stage decision task with reaction-times and drift-diffusion modeling. *PLoS Comput Biol* **15**:1–25. doi:10.1371/journal.pcbi.1006803

Spektor MS, Kellen D. 2018. The relative merit of empirical priors in non-identifiable and sloppy models: Applications to models of learning and decision-making. *Psychon Bull Rev* **25**:2047–2068. doi:10.3758/s13423-018-1446-5

Spiegelhalter DJ, Best NG, Carlin BP, Van der Linde A. 2002. Bayesian measures of model complexity and fit. *J R Stat Soc Ser B (Statistical Methodol* **64**:583–639.

Sutton, Richard S. 1988. Learning to Predict by the Method of Temporal Differences. *Mach Learn* **3**:9–44. doi:10.1023/A:1018056104778

Sutton RS, Barto AG. 2018. Reinforcement Learning: An Introduction, 2nd ed, MIT Press. Cambridge, MA: MIT press.

Teodorescu AR, Moran R, Usher M. 2016. Absolutely relative or relatively absolute: violations of value invariance in human decision making. *Psychon Bull Rev* **23**:22–38. doi:10.3758/s13423-015-0858-8

Ter Braak CJF. 2006. A Markov Chain Monte Carlo version of the genetic algorithm Differential Evolution: easy Bayesian computing for real parameter spaces. *Stat Comput* **16**:239–249. doi:10.1007/s11222-006-8769-1

Thura D, Cisek P. 2016. Modulation of Premotor and Primary Motor Cortical Activity during Volitional Adjustments of Speed-Accuracy Trade-Offs. *J Neurosci* **36**:938–956. doi:10.1523/JNEUROSCI.2230-15.2016

Tillman G, Van Zandt T, Logan GD. 2020. Sequential sampling models without random between-trial variability: the racing diffusion model of speeded decision making. *Psychon Bull Rev*. doi:10.3758/s13423-020-01719-6

940 Trueblood JS, Heathcote A, Evans NJ, Holmes WR. 2020. Urgency, Leakage, and the Relative Nature of
941      Information Processing in Decision-making. *Psychol Rev* 706291. doi:10.1101/706291
942 Turner BM. 2019. Toward a Common Representational Framework for Adaptation. *Psychol Rev*.
943      doi:10.1037/rev0000148
944 Turner BM, Sederberg PB, Brown SD, Steyvers M. 2013. A method for efficiently sampling from distributions
945      with correlated dimensions. *Psychol Methods* **18**:368–384. doi:10.1037/a0032222
946 van Maanen L, Miletić S. 2020. The interpretation of behavior-model correlations in unidentified cognitive
947      models. *Psychon Bull Rev*. doi:10.3758/s13423-020-01783-y
948 van Maanen L, van der Mijn R, van Beurden MHPH, Roijendijk LMM, Kingma BRM, Miletić S, van Rijn H.
949      2019. Core body temperature speeds up temporal processing and choice behavior under deadlines. *Sci Rep*
950      **9**:10053. doi:10.1038/s41598-019-46073-3
951 Van Ravenzwaaij D, Brown SD, Marley AAJ, Heathcote A. 2020. Accumulating advantages: A new
952      conceptualization of rapid multiple choice. *Psychol Rev* **127**:186–215. doi:10.1037/rev0000166
953 Voss A, Nagler M, Lerche V. 2013. Diffusion models in experimental psychology: A practical introduction. *Exp*
954      *Psychol* **60**:385–402. doi:10.1027/1618-3169/a000218
955 Voss A, Rothermund K, Voss J. 2004. Interpreting the parameters of the diffusion model: An empirical
956      validation. *Mem Cognit* **32**:1206–1220. doi:10.3758/BF03196893
957