

FRONT MATTER

Title

The other side of the coin: the paradoxical consequences of range adaptation in human reinforcement learning

Authors

Sophie Bavard^{1,2,3}, Aldo Rustichini⁴, Stefano Palminteri^{1,2,3}

Affiliations

1 Laboratoire de Neurosciences Cognitives et Computationnelles, Institut National de la Santé et Recherche Médicale, 29 rue d'Ulm 75005 Paris, FR.

2 Ecole normale supérieure, 29 rue d'Ulm 75005 Paris, FR.

3 Université de Recherche Paris Sciences et Lettres, 60 rue Mazarine 75006 Paris, FR.

4 University of Minnesota, 1925 4th Street South 4-101, Hanson Hall, Minneapolis, USA.

Abstract (150)

Evidence suggests that economic values are rescaled as a function of the range of the available options. Critically, although locally adaptive, range adaptation has been shown to lead to suboptimal choices. This is particularly striking in reinforcement learning (RL) situations when options are extrapolated from their original context. Range adaptation can be seen as the result of an adaptive coding process aiming at increasing the signal-to-noise ratio. However, this hypothesis leads to a counterintuitive prediction: decreasing task difficulty should increase range adaptation and, consequently, extrapolation errors. Here, we tested the paradoxical relation between range adaptation and performance in a large sample of participants performing variants of a RL task, where we manipulated task difficulty. Results confirmed that range adaptation induces systematic extrapolation errors and is stronger when decreasing task difficulty. Finally, we propose a range-adapting model and show that it is able to parsimoniously capture all the behavioral results.

MAIN TEXT

Introduction

In the famous Ebbinghaus illusion, two circles of identical size are placed near to each other, and larger circles surround one, while smaller circles surround the other. As a result, the central circle surrounded by larger circles appears smaller than the central circle surrounded by smaller circles, indicating that the subjective estimation of size of an object is affected by its the surroundings.

Beyond perceptual decision-making, wealth of evidence in neuroscience and in economics suggests that the subjective economic value of one option is not estimated in isolation, but is highly dependent of the context in which the options are presented (1,2). The vast majority of neuroeconomic studies of context-dependent valuation in humans considered situations where subjective values are triggered by explicit cues, that is stimuli whose value can be directly inferred, such as lotteries or snacks (3–5). However, in a series of recent papers, we and other groups demonstrated that contextual adjustments also permeate reinforcement learning situations, i.e., when option values have to be inferred from the history of past outcomes (6–8). We showed that an option, whose small objective value (7.5c) is learned in a context of smaller outcomes, is preferred to an option whose objective value (25c) is learned in a context of bigger outcomes, thus providing an economic equivalent of the Ebbinghaus illusion. Similar observations in birds suggest that this is a feature of decision-making broadly shared across vertebrates (9,10).

Although (as illustrated in the previous example) value context-dependence may lead to suboptimal decisions, it could be normatively understood as an adaptive process aimed at rescaling the behavioral response as a function of the range of the available options. Specifically, it could be seen as the result of an adaptive coding process aiming at increasing the signal-to-noise ratio by a system (the brain) constrained by the fact that behavioral variables have to be encoded by finite firing rates. In other terms, such *range adaptation* would be a consequence of how the system adjusts and optimizes the function associating the firing rate to the objective value to put its slope its maximum for each context (11,12).

If range adaptation is an automatic consequence of how the brain adapts its response to the distributions of the available outcomes, factors that facilitate the identifiability of these distributions should make it more pronounced. This would translate into a bigger difference between the objective option values (context-independent or *absolute*) and their corresponding subjective values (context-dependent or *relative*).

This leads to a counterintuitive prediction in the context of reinforcement learning. In fact, this is in striking contrast with the intuition embedded in virtually all learning algorithms, that making a learning-problem easier (by facilitating the identification of the outcome distributions) should, if anything, lead to more accurate and objective internal representations. In the present study, we aim at testing this hypothesis, while concomitantly gaining a better understanding of range adaptation at the computational level.

To empirically test this hypothesis, we build on previous research, and used a task featuring a *learning phase* and a *transfer phase* (6). In the learning phase, participants had to determine by trial-and-error the most favorable option in four fixed pairs of options (contexts), with different outcome ranges. In the transfer phase, the original options were rearranged, thus creating new contexts. This setup allowed us to quantify learning (or acquisition) errors during the first phase, and transfer (or extrapolation) errors during the second phase. Crucially, the task contexts were designed such that the correct responses in the transfer phase presented an overall higher expected value. We varied this paradigm in eight different versions where we manipulated the task difficulty in complementary ways. First, some of the experiments (E3, E4, E7, E8) featured *complete* feedback information, meaning that participants were informed about the outcome of the forgone option. This manipulation reduces task difficulty by resolving the uncertainty

concerning the counterfactual outcome. Accordingly, it has been repeatedly shown to improve learning performance (8,13). Second, some of the experiments (E5, E6, E7, E8) featured a block (instead of interleaved design), meaning that all the trials featuring one context were presented in a row. This manipulation reduces task difficulty by reducing working memory demand and has also been shown to improve learning performance (14). Finally, in some of the experiments (E2, E4, E6, E8), feedback was also provided in the transfer phase, thus allowing to assess if and how the values learned during the learning phase can be revised.

Behavioral analyses backed up our prediction and indicate that acquisition error rate in the learning phase is largely dissociable from extrapolation error rate in the transfer phase. Critically (and paradoxically), transfer phase error rate was higher when the learning phase was easier. Accordingly, the estimated deviation between the objective values and the subjective values increased in the complete feedback and block design tasks. The deviation was corrected only in the experiments featuring complete feedback in the transfer test.

To complement choice rate analysis, we developed a computational model that implements range adaption as a range normalization process, by tracking the maximum and the minimum possible reward in each learning context. Model simulations parsimoniously captured performance in the learning and the transfer phase, including the range adaptation-induced suboptimal preferences. Model simulations also allowed us to rule out alternative interpretations of our results that could come from two prominent psychological and economic theories: habit formation and risk aversion (15,16). Model comparison results were confirmed by checking out-of-sample likelihood as a quantitative measure of goodness of fit.

1 **Results**2 **Experimental protocol**

3 We designed a series of learning and decision-making experiments involving variants of a
4 main task. The main task was composed of two phases: the *learning* and the *transfer* phase. During
5 the learning phase, participants were presented with eight abstract pictures, organized in four stable
6 choice contexts. In the learning phase, each choice context featured only two possible outcomes:
7 either 10pt/0pt or 1pt/0pt. The outcomes were probabilistic (75% or 25%) and we labeled the
8 choices contexts as a function of the difference in expected value between the most and the least
9 rewarding option: $\Delta EV=5$ and $\Delta EV=0.5$ (**Figure 1A**). In the subsequent transfer phase, the eight
10 options were re-arranged into new choice contexts, where options associated with 10pt were
11 compared to options associated with 1pt (see (7,10) for similar designs in humans and starlings).
12 The resulting new four contexts were labeled $\Delta EV=7.25$, $\Delta EV=6.75$, $\Delta EV=2.25$, $\Delta EV=1.75$
13 (**Figure 1B**). In our between-subjects study, we developed eight different variants of the main
14 paradigm where we manipulated whether we provided trial-by-trial feedback in the transfer phase
15 (with / without), the quantity of information provided at feedback (partial: only the outcome of the
16 chosen option is shown / complete: both outcomes are shown) and the temporal structure of choice
17 contexts presentation (interleaved: choice contexts appear in a randomized order / block: all trials
18 belonging to the same choice contexts are presented in a row) (**Figure 1C**). All the experiments
19 reported in the main text were conducted online (N=100 participants in each version of the
20 experiment); we report in the supplementary information the results concerning a similar
21 experiment realized in the lab (see **Supplementary Materials**).

22

23 **Overall correct response rate**

24 The main dependent variable in our study was the correct response rate, i.e., the proportion of
25 expected value-maximizing choices in the learning and the transfer phase (crucially our task design
26 allowed to identify an expected value-maximizing choice in all choice contexts). In the learning
27 phase, the average correct response rate was significantly higher than chance level 0.5 (0.69 ± 0.16 ,
28 $t(799) = 32.49$, $p < .0001$, $d = 1.15$; **Figure 2A-B**). Replicating previous findings, in the learning
29 phase, we also observed a moderate but significant effect of the choice contexts, where the correct
30 choice rate was higher in the $\Delta EV=5.0$ compared to the $\Delta EV=0.5$ contexts (0.71 ± 0.18 vs
31 0.67 ± 0.18 ; $t(799) = 6.81$, $p < .0001$, $d = 0.24$; **Figure 2C**)(6).

32 Correct response rate was also higher than chance in the transfer phase (0.62 ± 0.17 , $t(799) = 20.29$,
33 $p < .0001$, $d = 0.72$, **Figure 2D-E**), but it was also strongly modulated by the choice context
34 ($F(2.84, 2250.66) = 271.68$, $p < .0001$, $\eta^2 = .20$, Huynh–Feldt corrected). In the transfer phase, the
35 $\Delta EV=1.75$ choice context is of particular interest, since the expected value maximizing option was
36 the least favorable option of a $\Delta EV=5.0$ context in the learning phase, and, conversely, the expected
37 value minimizing option was the most favorable option of a $\Delta EV=0.5$ context of the learning phase.
38 In other words, for subject relying on expected values calculated in a context-independent scale,
39 the $EV_{2.5}$ option is preferred compared to the $EV_{0.75}$ option. On the other side, a subject encoding
40 the option values on a fully context-independent manner (which is equivalent to encode the rank
41 between two options in a given context), will perceive the $EV_{2.5}$ option as the less favorable
42 option compared to the $EV_{0.75}$. Therefore, preferences in the $\Delta EV=1.75$ context are diagnostic of

43 whether values are learned and encoded in an absolute or relative scale. Crucially, in the $\Delta EV=1.75$
44 context, we found that participants' average correct choice rate was significantly *below* chance
45 level (0.42 ± 0.30 , $t(799) = -7.25$, $p < .0001$, $d = -0.26$; **Figure 2F**), thus demonstrating that
46 participants express suboptimal preferences in this context, i.e., they do not choose the option with
47 the highest objective expected value.

48

49 **Between-experiments comparisons: learning phase**

50 In this section we analyze the correct response rate as a function of the experimental factors
51 manipulated across the eight experiments (the quantity of provided information, that could be either
52 partial or complete; the temporal structure of choice contexts presentation, that could be block or
53 interleaved; and whether feedback was provided in the transfer phase). In the main text we report
54 the significant results, but please see **Tables 1 and 2** for all results and effect sizes.

55 First, we analyzed the correct choice rate in the learning phase (**Figure 2B**). As expected,
56 increasing feedback information had a significant effect on correct choice rate in the learning phase
57 ($F(1,792) = 55.57$, $p < .0001$, $\eta_p^2 = .18$); similarly, performance in the block design experiments
58 was significantly higher ($F(1,792) = 87.22$, $p < .0001$, $\eta_p^2 = .25$). We found a significant interaction
59 between feedback information and task structure, reflecting that the difference of performance
60 between partial and complete feedback was higher in block design ($F(1,792) = 5.05$, $p = .02$, $\eta_p^2 =$
61 $.02$). We found no other significant main effect, double or triple interaction (**Table 1**).

62 We also analyzed the difference in performance between the $\Delta EV=5.0$ and $\Delta EV=0.5$ choice
63 contexts across experiments (**Figure 2C**). We found a small but significant effect of temporal
64 structure, the differential being smaller in the block compared to interleaved experiments ($F(1,792)$
65 $= 7.71$, $p = .006$, $\eta_p^2 = .01$), and found no other significant main effect, nor interaction.

66 To sum up, as expected (8,13,14), increasing feedback information and clustering the choice
67 contexts had a beneficial effect on correct response rate in the learning phase. Designing the choice
68 contexts in blocks also blunted the difference in performance between the small ($\Delta EV=0.5$) and
69 big ($\Delta EV=5.0$) magnitude contexts.

70

71 **Between-experiments comparisons: transfer phase**

72 We then analyzed the correct choice rate in the transfer phase (**Figure 2E**). Unsurprisingly,
73 showing trial-by-trial feedback in the transfer phase led to significantly higher performance
74 ($F(1,792) = 137.18$, $p < .0001$, $\eta_p^2 = .07$). Increasing feedback information from partial to complete
75 also had a significant effect on transfer phase correct choice rate ($F(1,792) = 22.36$, $p < .0001$, η_p^2
76 $= .01$). Interestingly, we found no significant main effect of task structure in the transfer phase (see
77 **Table 1**).

78 We found a significant interaction between feedback information and the presence of feedback in
79 the transfer phase, showing that the increase in performance due to the addition of feedback
80 information is higher when both outcomes were displayed during the learning phase ($F(1,792) =$
81 20.18 , $p < .0001$, $\eta_p^2 = .01$). We also found a significant interaction between transfer feedback and
82 task structure, reflecting that the increase in performance due to the addition of feedback
83 information was even higher in block design ($F(1,792) = 42.22$, $p < .0001$, $\eta_p^2 = .02$). Finally, we
84 found a significant triple interaction between feedback information, the presence of feedback in the

85 transfer phase, and task structure ($F(1,792) = 5.02, p = .03, \eta_p^2 = .003$). We found no other
86 significant double interaction. We also separately analyzed the correct choice rate in the $\Delta EV=1.75$
87 context (**Figure 2F**). Overall, the statistical effects presented a similar pattern as the correct choice
88 rate across all conditions (see **Table 2**), indicating that overall correct choice rate and the correct
89 choice rate in the key comparison $\Delta EV=1.75$ provided a coherent picture. Furthermore, comparing
90 the $\Delta EV=1.75$ to chance level (0.5) revealed that participants, overall, significantly expressed
91 *reward minimizing preferences* in this choice context. Crucially, the lowest correct choice rate was
92 observed in the experiment featuring complete feedback, clustered choice contexts and no feedback
93 in the transfer phase (E7; $0.27 \pm 0.32, t(99) = -7.11, p < .0001, d = -0.71$); the addition of feedback
94 in the transfer phase reversed the situation, since the only experiment where participants expressed
95 reward maximizing preference was E8 ($0.59 \pm 0.29, t(99) = 2.96, p = .0038, d = 0.30$).

96

97 **Between-phase comparison**

98 Interestingly, we found a significant interaction between the phase (learning or transfer) and
99 transfer feedback (without/with) on correct choice rate ($F(1,792) = 82.30, p < .0001, \eta_p^2 = .09$).
100 This interaction is shown in **Figure 3** and reflects the fact that while adding transfer feedback
101 information had a significant effect on transfer performance ($F(1,792) = 137.18, p < .0001, \eta_p^2 =$
102 $.05, \text{Figure 3A-B}$), it was not sufficient to outperform learning performance (with transfer
103 feedback: learning performance 0.69 ± 0.16 vs transfer performance $0.68 \pm 0.15, t(399) = 0.89, p =$
104 $.38, d = 0.04, \text{Figure 3B}$).

105 Finally, close inspection of the learning curves revealed that in experiments where feedback was
106 not provided in the transfer phase (E1, E3, E5 and E7), choice rates (and therefore option
107 preferences) were stationary (**Figure 3A** and **Figure 3B**). This observation rules out the possibility
108 that reduced performance in the transfer phase was induced by progressively forgetting the values
109 of the options (in which case we should have observed a non-stationary and decreasing correct
110 response rate).

111 In conclusion, comparison between the learning and the transfer phase reveals two inter-related
112 and intriguing facts: i) despite the fact that the transfer phase happens immediately after an
113 extensive learning phase, performance is, if anything, lower compared to the learning phase; ii)
114 factors that improve performance (by intrinsically or extrinsically reducing task difficulty) in the
115 learning phase have either no (feedback information) or a negative (task structure) impact on the
116 transfer phase performance.

117

118 **Inferred option values**

119 To visualize and quantify how much observed choices deviate from the experimentally determined
120 true option values, we optimized the four possible option values as free parameters.

121 More precisely, we initialized each subjective value at their true value (we labeled the four possible
122 expected values as follows: $EV_{7.5}, EV_{2.5}, EV_{0.75}$, and $EV_{0.25}$), and optimized these values by
123 gradient descent in order to maximize the likelihood of observing participants' choices using the
124 logistic function (for, say, options $EV_{2.5}$ and $EV_{0.75}$):

125

$$P(EV_{2.5}) = \frac{1}{1 + e^{(V(EV_{0.75}) - V(EV_{2.5}))}} \quad (1)$$

126 So that, if a participant chose indifferently between the $EV_{2.5}$ and the $EV_{0.75}$ option, their fitted
127 values would be very similar: $V(EV_{2.5}) \approx V(EV_{0.75})$. Conversely, a participant with a sharp
128 (optimal) preference for $EV_{2.5}$ over $EV_{0.75}$ would lead different fitted values: $V(EV_{2.5}) > V(EV_{0.75})$.
129 In a first step, in the experiments where feedback was not provided in the transfer phase (E1, E3,
130 E5 and E7), we optimized a set of subjective values per participant.

131 Consistent with the correct choice rate results described above, we found a value inversion of the
132 two intermediary options ($EV_{2.5}$ 4.46 ± 1.2 , $EV_{0.75}$ 5.26 ± 1.2 , $t(399) = -7.82$, $p < .0001$, $d = -0.67$),
133 which were paired in the $\Delta EV = 1.75$ context (**Figure 3C**). The differential was also strongly
134 modulated across experiments ($F(3,396) = 18.9$, $p < .0001$, $\eta_p^2 = .13$, **Figure 3C**) and reached its
135 highest value in E7 (complete feedback and block design).

136 As a second step, in the experiments where feedback was provided in the transfer phase (E2, E4,
137 E6 and E8), we optimized a set of subjective values per trial. This fit allows us to estimate the trial-
138 by-trial evolution of the subjective values over task time. The results of this analysis clearly show
139 that suboptimal preferences progressively arise during the learning phase and disappear during the
140 transfer phase (**Figure 3D**). However, the suboptimal preference was completely corrected only in
141 E8 (complete feedback, block design) by the end of the transfer phase.

142 The analysis of inferred option values clearly confirms that participants' choices do not follow the
143 true underlying objective monotonic ordering of the option values. Furthermore, it also clearly
144 illustrates that in choice contexts that are supposed to facilitate the learning of the option values
145 (complete feedback, block design), the deviation from monotonic ordering, at least at the beginning
146 of transfer phase, is paradoxically greater. Monotonicity was fully restored only in E8, where
147 complete feedback was provided in the transfer phase.

148

149 **Computational formalization of the behavioral results**

150 To formalize context-dependent reinforcement learning and account for the behavioral results, we
151 designed a modified version of a standard model, where option-dependent Q-values are learnt from
152 a range-adapted reward term. In the present study we implemented range adaptation as a range
153 normalization process, which one among other possible implementations (17). At each trial t , the
154 relative reward, $R_{RAN,t}$, is calculated as follows:

$$155 \quad R_{RAN,t} = \frac{R_{ABS,t} - R_{MIN,t}(s)}{R_{MAX,t}(s) - R_{MIN,t}(s) + 1} \quad (2)$$

156 where s is the decision context (i.e., a combination of options) and R_{MAX} is a state-level variable,
157 initialized to 0 and updated at each trial t if the outcome is greater than its current value. As such,
158 R_{MAX} will converge to the maximum outcome value in each decision context, which in our task is
159 either 1pt or 10pt. In the first trial $R_{RAN} = R_{ABS}$ (because $R_{MAX,0}(s) = 0$), and in later trials it is
160 progressively normalized between 0 and 1 as the range value $R_{MAX}(s)$ converges to its true value.
161 Since in our task the minimum possible outcome is always zero $R_{MIN,t}$ update was omitted while
162 fitting the main experiments (but included in a ninth dataset analyzed below).

163 We refer to this model as the RANGE model and we compared it to a benchmark model
164 (ABSOLUTE) which updates option values based the absolute reward values (note that the
165 ABSOLUTE is nested within the RANGE model).

166 For each model, we estimated the optimal free parameters by likelihood maximization. We used
167 the out-of-sample likelihood to compare goodness-of-fit and parsimony of the different models. To
168 calculate the out-of-sample likelihood in the learning phase, the optimization was performed on
169 half of the trials (one $\Delta EV=5.0$ and one $\Delta EV=0.5$ decision context) in the learning phase, and the
170 best fitting parameters in this first set were used to predict choices in the remaining half of trials.
171 In the learning phase, we found that the RANGE model significantly outperformed the
172 ABSOLUTE model (out-of-sample LL_{RAN} vs. LL_{ABS} , $t(799) = 6.89$, $p < .0001$, $d = 0.24$, **Table 3**).
173 To calculate the out-of-sample likelihood in the transfer phase, the optimization was performed on
174 all trials of the learning phase and the best fitting parameters in the learning phase were used to
175 predict choices in the transfer phase. Thus, the resulting likelihood is not only out-of-sample, but
176 also cross-learning phase. This analysis revealed that the RANGE model outperformed the
177 ABSOLUTE model (out-of-sample LL_{RAN} vs. LL_{ABS} , $t(799) = 8.56$, $p < .0001$, $d = 0.30$).

178 To study the behaviors of our computational model and confirm the behavioral reasons underlying
179 the out-of-sample likelihood results, we simulated the two models (using the individual best fitting
180 parameters)(18). In the learning phase, only the RANGE model managed to reproduce the observed
181 correct choice rate. Specifically, the ABSOLUTE model predicts very poor performance in the
182 $\Delta EV=0.5$ context (ABS vs. data, $t(799) = -16.90$, $p < .0001$, $d = 0.60$, RAN vs. data, $t(799) = -1.79$,
183 $p = .07$, $d = -0.06$, **Figure 4A**).

184 In the transfer phase, and particularly in the $\Delta EV=1.75$ context, only the RANGE model manages
185 to account for the observed correct choice rate, while the ABSOLUTE model fails (ABS vs. data
186 $t(799) = 13.20$, $p < .0001$, $d = 0.47$, RAN vs. data $t(799) = 0.36$, $p = .72$, $d = 0.01$, **Figure 4C-D**).
187 In general, the ABSOLUTE model tends to overestimate the correct choice rate in the transfer
188 phase.

189 In addition to looking at the choice patterns, we submitted the RANGE model simulations to the
190 inferred subjective option values analysis. Not only is the RANGE model able to capture the value
191 inversion that we observed in the data, as well as the estimated options (RAN vs data, $t(799) =$
192 1.55 , $p = .12$, $d = 0.06$, **Figure 4E**), but it is also able to predict its dynamic emergence and its trial-
193 by-trial evolution (**Figure 4F**).

194

195 **Ruling out habit formation**

196 One of the distinguishing behavioral signatures of the RANGE model compared the ABSOLUTE
197 one is the preference for the suboptimal option in the $\Delta EV=1.75$ context. Since the optimal option
198 in the $\Delta EV=1.75$ context is not often chosen during the learning phase (where it is locally
199 suboptimal), it could be argued that this result arises from taking decisions based on a weighted
200 average between their absolute values and past choice propensity (a sort of *habituation* or choice
201 trace). To rule out this interpretation, we fitted and simulated a version of a HABIT model, which
202 takes decisions based on a weighted sum of the absolute Q-values and a habitual choice trace
203 (16,19). The habitual choice trace component is updated with an additional learning rate parameter
204 that gives a bonus to the selected action. Decisions are taken comparing option-specific decision-
205 weights D_t :

$$206 \quad D_t(s, c) = (1 - \omega) * Q_t(s, c) + \omega * H_t(s, c) \quad (3)$$

207 where at each trial t , state s and chosen option c , ω is the arbiter, Q is the absolute Q-value H is the
208 habitual choice trace component. The weight ω is fitted as an additional parameter (for $\omega=0$ the
209 model reduces to the ABSOLUTE model) and governs the relative influence of each controller.

210 We found that the HABIT model, similarly to the ABSOLUTE model, fails to perfectly match the
211 participants' behavior, especially in the $\Delta EV=0.5$ and $\Delta EV=1.75$ contexts (**Figure 5A**). Indeed, in
212 the learning phase, the addition of a habitual component is not enough to cope for the difference in
213 option values, and therefore the model simulations in the transfer phase fail to match the observed
214 choice pattern (**Figure 5B**). This is because the HABIT model encodes values on an absolute scale
215 and does not manage to develop a strong preference for the correct response in the $\Delta EV=0.5$
216 context, in the first place (**Figure 5A**). Thus, it does not carry a choice trace strong enough to
217 overcome the absolute value of the correct response in the $\Delta EV=1.75$ context (**Figure 5B, Supp.**
218 **Figure 2 A-B, Table 3**).

219 To summarize, a model assuming absolute value encoding coupled with a habitual component
220 could not fully explain observed choices in both the learning and transfer phase.

221

222 **Ruling out diminishing marginal utility**

223 One of the distinguishing behavioral signatures of the RANGE model is that it predicts very similar
224 correct choice rates in the $\Delta EV=5.00$ and the $\Delta EV=0.50$ contexts compared to the behavioral data,
225 while both the ABSOLUTE and the HABIT predict a huge drop in performance in the $\Delta EV=0.50$
226 that directly stems from its small difference in expected value. It could be argued that this result
227 arises from the fact that expected *utilities* (and not expected *values*) are learned in our task.
228 Specifically, a diminishing marginal utility parameter would blunt differences in outcome
229 magnitudes and would suppose that choices are made by comparing outcome probabilities. The
230 process could also explain the preference for the suboptimal option in the $\Delta EV=1.75$ context, since
231 the optimal option in the $\Delta EV=1.75$ context is rewarded (10pt) only the 25% of the time, while the
232 suboptimal option is rewarded (1pt) 75% of the time. To rule out this interpretation, we fitted and
233 simulated a UTILITY model, which updates Q-value based reward utilities calculated from
234 absolute reward as follows:

$$235 \quad R_{UTI,t} = (R_{ABS,t})^{\nu} \quad (4)$$

236 where the exponent ν is the utility parameter ($0 < \nu < 1$, for $\nu = 1$ the model reduces to the
237 ABSOLUTE model).

238 We found that the UTILITY model, similarly to the RANGE model, captures quite well the
239 participants' behavior in the $\Delta EV=0.5$ context (**Figure 5C**). However, concerning the transfer
240 phase (especially the $\Delta EV=1.75$ context), it fails to capture the observed pattern (**Figure 5C-D**).
241 Additional analyses suggest that this is specifically driven by the experiments where the feedback
242 was provided during the transfer phase (**Figure 5D**). Indeed, the static nature of the UTILITY fails
243 to match the fact that the preferences in the $\Delta EV=1.75$ context can be reversed by providing
244 complete feedback (**Supp. Figure 2 C-D**). Model comparison showed that the RANGE model also
245 outperformed the UTILITY model (out-of-sample LL_{RAN} vs. LL_{UTY} , $t(799) = 3.21$, $p = .001$, $d =$
246 0.06 , **Table 3**). The comparison between the RANGE and the UTILITY model went in the same
247 direction with a similar magnitude, reaching marginal statistical significance (out-of-sample LL_{RAN}
248 vs. LL_{HAB} , $t(799) = 1.77$, $p = .07$, $d = 0.05$). To summarize, a model assuming diminishing marginal
249 utilities could not fully explain observed choices in the transfer phase.

250

251 **Sub-optimality of range-adaptation in our task**

252

253 The RANGE model is computationally more complex compared to the ABSOLUTE model, as it
254 presents an additional internal variable (R_{MAX}), which is learnt with a dedicated parameter. Here
255 we wanted to assess whether this additional computational complexity really paid off in our task.

256 We split the participants according to the sign of out-of-sample likelihood difference between the
257 RANGE and the ABSOLUTE model: if positive, the RANGE model better explains the
258 participant's data (RAN>ABS), if negative, the ABSOLUTE model does (ABS>RAN). Reflecting
259 our overall model comparison result, we found more participants in the RAN>ABS, compared to
260 the ABS>RAN category (N=545 vs. N=255).

261 We found no main effect of winning model on overall (both phases) performance ($F(1,798) = 0.03$,
262 $p = .87$, $\eta_p^2 = 0$). Interestingly, we found that while RANGE encoding is beneficial and allows for
263 better performances in the learning phase, it leads to worst performance in the transfer phase
264 ($F(1,798) = 187.3$, $p < .0001$, $\eta_p^2 = .19$, **Figure 6A**). In other terms, in our task, it seems that the
265 learning phase and the transfer phase are playing the game tug of war: when performance are pulled
266 in favor of the learning phase (RANGE model participants) this will be at the cost of the transfer
267 phase (and vice versa).

268 A second question is whether overall in our study, behaving as a RANGE model revealed
269 economically advantageous. To answer this question, we compared the final monetary payoff in
270 the real data, following the simulations using the participant-level best fitting parameters.
271 Consistently with the task design, we found that the monetary outcome was higher in the transfer
272 phase than in the learning phase (transfer gains $M = 2.16 \pm 0.54$, learning gains $M = 1.99 \pm 0.35$,
273 $t(799) = 8.71$, $p < .0001$, $d = 0.31$). Crucially, we found that the simulation of the RANGE model
274 induces significantly lower monetary earnings (ABS vs RAN, $t(799) = 19.39$, $p < .0001$, $d = 0.69$,
275 **Figure 6B**). This result indicates that despite being locally adaptive (in the learning phase), in our
276 task, range adaptation is economically disadvantageous, thus supporting the idea that it is the
277 consequence of automatic, uncontrolled, process.

278

279 **Validation of range adaptation in previous dataset**

280

281 Our main experiments only featured positive outcomes, in addition to 0. Since in our model the
282 state-level variables (R_{MAX} and R_{MIN}) are initialized at 0, R_{MAX} converges to the maximum
283 outcome value in each choice context, while R_{MIN} remains 0 in every trial and choice context. This
284 is set-up is not ideal to test the full normalization rule we are proposing here. To obviate this
285 limitation, we re-analyzed a ninth dataset (N=60) from a previously published study on a related
286 topic (6). Crucially, in addition to an outcome magnitude manipulation ('10c vs '1€', similar to our
287 learning phase), this study also manipulated the valence of the outcomes (gain vs loss). This latter
288 manipulation allows to assess situations where the value of R_{MIN} can change and converge to a
289 negative values, thus allowing us to compare the full range normalization rule to a simplified
290 version:

291

292

$$\frac{R_{ABS} - R_{MIN}}{R_{MAX} - R_{MIN}} \text{ vs } \frac{R_{ABS}}{R_{MAX}}$$

293 We later refer to the simplified version of the model as the RMAX model. Model simulations show
294 that while the RMAX model can capture the learning and transfer patterns for the gain-related
295 options, it fails to do so for the loss-related options (**Figure 5E-F**). Indeed, the lack of update of
296 R_{MAX} in the loss contexts induces the RMAX model to encode values on an absolute scale (without
297 normalization) and therefore, to fail to account for the range adaptation process. On the other hand,
298 by updating both R_{MAX} in the gain contexts and R_{MIN} in the loss contexts, the RANGE model
299 adapts with a full range in all contexts and is able to match participants' behavior even in the loss-
300 related options as well as the gain-related options, in both the learning and transfer phases (**Figure**
301 **5E-F**). To conclude, this final analysis provides crucial support to the idea that range adaptation is
302 consistent with a full range normalization rule, which takes into account both the maximum and
303 the minimum possible outcomes.

304
305
306
307

Discussion

308 In the present paper we investigate context-dependent reinforcement learning, more specifically
309 range adaptation, in a large cohort of human participants tested online over eight different variants
310 of a behavioral task. Building on previous studies of context-dependent learning, the core idea of
311 the task is to juxtapose an initial learning phase with fixed pairs of options (featuring either small
312 or big outcomes) to a subsequent transfer phase where options are rearranged in new pairs (mixing
313 up small and big outcomes)(6,7,10). In some experiments, we directly reduced task difficulty by
314 reducing outcome uncertainty providing complete feedback. In some experiments, we indirectly
315 modulated task difficulty by clustering in time the trials of a given contexts; therefore, reducing
316 working memory demand. Finally, in some experiments, feedback was also provided in the transfer
317 phase.

Behavioral findings

318 As expected, correct choice rate in the learning phase was higher when the feedback was complete,
319 which indicates that participants integrated the outcome of the forgone option when it is presented
320 (8,14). Also expectedly, in the learning phase participants displayed a higher correct choice rate
321 when the trials of a given context were all blocked together, indicating that reducing working
322 memory demands facilitate learning (15). Replicating previous findings, we also found that,
323 overall, correct response rate was slightly but significantly higher in the big magnitude contexts
324 ($\Delta EV=5.0$), but the difference was much smaller compared to what one would expect assuming
325 absolute value learning and representation (as showed by the ABSOLUTE model simulations (6)):
326 a pattern consistent with a *partial* range adaptation. The outcome magnitude-induced difference in
327 correct choice rate was significantly smaller and not different from zero in block experiments (*full*
328 adaptation), thus providing a first hint that reducing task difficulty increases range adaptation.
329 Despite learning phase performance being fully consistent with our hypothesis, the crucial evidence
330 comes from the results of the transfer phase. First, overall correct response rate pattern in the
331 transfer phase did not follow that of the learning phase. Complete feedback and block design factors
332 have no direct beneficial effect on transfer phase performance. In fact, if anything, the *worst*
333 possible transfer phase performance was obtained in a complete feedback and block experiment.
334 This was particularly striking in the $\Delta EV=1.75$ condition, where participants significantly preferred
335 the suboptimal option and, again, worst score was obtained in a complete feedback and block
336 design experiment. Second, we ruled out that the comparably low performance in the transfer phase
337

338 was due to having forgotten the value of the options. Indeed, since the transfer phase is, by
339 definition, after the learning phase, although very unlikely (the two phases were only few seconds
340 apart), it is conceivable that a drop in performance is due to the progressive forgetting of the option
341 values. Two features of the correct choice rate curves allowed to reject this interpretation: i) correct
342 choice rate abruptly decreases just after the learning phase; ii) when feedback is not provided the
343 choice rate remains perfectly stable with no sign of regression to chance level. On the other side,
344 i.e., when feedback was provided in the transfer phase, the correct choice rate increased to reach
345 (on average) the level reached at the end of the learning phase. The results are therefore consistent
346 with the idea that in the transfer phase, participants express context-dependent option values
347 acquired during the learning phase, which entails a first counterintuitive phenomenon: even if the
348 transfer phase is performed immediately after the learning phase, the correct choice rate drops. This
349 is due to the rearrangement of the options in new choice contexts, where options that were
350 previously optimal solutions (in the small magnitude contexts) become suboptimal solutions. We
351 also observed a second counterintuitive phenomenon: factors that increase performance during the
352 learning phase (i.e., increasing feedback information and reducing working memory load),
353 paradoxically further undermined transfer phase correct choice rate. The conclusions based on
354 these behavioral observations were confirmed by inferring the most plausible option values based
355 on the observed choices, where we could compare the objective ranking of the options to their
356 subjective estimation. The only experiment where we observed an almost monotonic ranking was
357 the partial feedback / interleaved experiment, even if we observed no significant difference between
358 the EV=2.5 and the EV=0.75 options. In all the other experiments, the EV=0.75 option was valued
359 more compared to the EV=2.5 option, with the highest difference observed in the complete
360 feedback / block design. Thus, in striking opposition with the almost universally shared intuition
361 that reducing task difficulty should lead to more accurate subjective estimates, here we present a
362 rare instance where the opposite is true.

363 **Computational mechanisms**

364 The observed behavioral results were satisfactorily captured by a parsimonious model (the RANGE
365 model) that instantiated a dynamic range normalization process. Specifically, the RANGE model
366 learns in parallel context-dependent variables (R_{MAX} and R_{MIN}) that are used to normalize the
367 outcomes. The R_{MAX} and R_{MIN} are learnt incrementally and the speed determines the extent of the
368 normalization, leading to partial or full range adaptation as a function of the contextual learning
369 rate. Developing a new model was necessary, as previous models of context-dependent
370 reinforcement learning did not include range adaptation and focused on different dimensions of
371 context-dependence (reference point-centering and outcome comparison) (7,8). The model also
372 represents an improvement over a previous study where we instantiated partial range adaptation
373 assuming a perfect and innate knowledge about the outcome ranges and a static hybridization
374 between relative and absolute outcome values (6).

375 One limitation is that in present formulation R_{MAX} and R_{MIN} can only grow. Again, this is a feature
376 that is well suited for our task, but may not correspond to many other laboratory-based and real-
377 life situations, where the range can drift over time. This limitation could be overcome by assuming,
378 for example, that R_{MAX} is also updated at a smaller rate when the observed outcome is smaller than
379 the current R_{MAX} (the opposite could be true R_{MIN}). Finally, we note that our model applied to the
380 main eight experiments (where R_{MIN} was irrelevant) can also be seen as special case of a divisive
381 normalization process (temporal normalization (20)). To verify the relevance of the full range
382 normalization rule, we re-analysed a previous dataset involving negative outcomes, where we were

383 able to show that both the R_{MAX} and R_{MIN} were important to explain the full spectrum of the
384 behavioral results. Finally, it is worth noting that range normalization has been shown to perform
385 poorly in explaining context-dependent decision-making in other (i.e., not reinforcement learning)
386 paradigms (17,21,22), opening to the possibility that the normalization algorithm is different in
387 experience-based and description-based choices. Future research contrasting different outcome
388 ranges, multiple-options tasks are required to firmly determine which the functional form of
389 normalization follows is better suited for experience-based and description-based choices (23).

390 We compared and ruled out another plausible computational interpretation derived from prominent
391 psychological theory (24,25). First, we considered a habit formation model (16). We reasoned that
392 our transfer phase results (and particularly the value inversion in the $\Delta EV=1.75$ context) could
393 derive from the participants choosing based on a weighted average between absolute values and
394 past choice propensities. In fact, the suboptimal option in the $\Delta EV=1.75$ context ($EV=0.75$) was
395 chosen more frequently than the optimal option ($EV=2.5$) in the learning phase. However, model
396 simulations showed that the HABIT model was not capable to explain the observed pattern. In fact,
397 in the learning phase, the HABIT model, just like the ABSOLUTE model, did not develop a
398 preference for the $EV=0.75$ option strong enough to generate a habitual trace sufficient to explain
399 the transfer phase pattern. Beyond model simulation comparisons, we believe that this
400 interpretation could have been rejected based on a priori arguments. The HABIT model can be
401 conceived as a way to model habitual behavior, i.e., responses automatically triggered by stimulus-
402 action associations. However, both in real life and laboratory experiments, habits have been shown
403 to be acquired over time scales (days, months, year) order of magnitudes bigger compared to the
404 timeframe of our tasks (26,27). Indeed, it is even debatable whether in our task participants
405 developed even a sense of familiarity toward the (never seen before) abstract cues we used as
406 stimuli. The HABIT model can also be conceived as a way to model choice hysteresis, sometimes
407 referred to as choice repetition of perseveration bias, that could arise from a form of sensory-motor
408 facilitation, where recently performance actions become facilitated (19,28). However, the screen
409 position of the stimuli was randomized in a trial-by-trial basis; most of the experiments involved
410 inter-leaved design, thus precluding any strong role for sensory-motor facilitation-induced choice
411 inertia.

412 We compared and ruled out a plausible computational interpretation derived from economic theory
413 (29). Since Bernoulli (1700-1782), risk aversion is explained by assuming diminishing marginal
414 utility of objective outcomes (30). At the limit, if diminishing marginal utility was applied in our
415 case, the utility of 10 points could be perceived as the utility of 1 point. In this extreme scenario,
416 choices would be only based on the comparison between the outcome probabilities. This could
417 explain most aspects of the choice pattern. Indeed, the UTILITY model did a much better job
418 compared to HABIT model. However, compared to the RANGE model, it failed to reproduce the
419 observed behavior of the experiments where feedback was provided in the transfer phase. This
420 naturally results from the fact that the model assumes diminishing marginal utility as being a static
421 property of the model and therefore cannot account for experience-dependent correction of context-
422 dependent biases. However, also in this case, a priori considerations could have ruled out the
423 UTILITY interpretation. Our experiment involves stakes small enough to make diminishing
424 marginal utility not reasonable. Rabin provides a full treatment of this issue, and shows that the
425 explaining risk aversion for small stakes (as those used in the laboratory) using diminishing
426 marginal utility leads to extremely unlikely prediction, such as turning down gambles with infinite
427 positive expected values (15). Indeed, following H. Markowitz intuition, most realistic models of
428 the utility function suppose risk neutrality (or risk seeking) for small gains (31).

429 Our results contribute to the old and still ongoing debate about whether the brain computes option-
430 oriented values independently from the decision-process itself (2,32). On one side of the spectrum,
431 decision theories such as expected utility theory and prospect theory, postulate that a value is
432 attached to each option independently of the other options simultaneously available (32). On the
433 other side of the spectrum, other theories, such as regret theory, postulate that the value of an option
434 is primarily determined by the comparison with other available options (33). A similar gradient
435 exists in the reinforcement learning framework, between methods such as the Q-learning, on one
436 side, and direct policy learning without value computations, on the other side (34). Recent studies
437 in humans, coupling imaging to behavioral modeling, provided some support for direct policy
438 learning in humans, by showing that, in complete feedback tasks, participants' learning was driven
439 by a teaching signal, essentially determined by the comparison between the obtained and the
440 forgone outcomes (regret/relief)(7,35). Beyond behavioral model comparison, analysis of neural
441 activity in the ventral striatum (a brain system traditionally thought to encode option-specific
442 prediction errors (36)), was also consistent with direct policy learning. However, while our findings
443 clearly falsify the Q-learning's assumption that option-values are learned in an absolute (or context-
444 independent) scale, model simulations also reject the other extreme view of direct policy learning
445 (see **Supplementary Materials**). Indeed, our results are rather consistent with a hybrid scenario
446 where option-specific values are initially encoded in an absolute scale and are progressively
447 normalized to eventually represent the context-specific rank of each option. This view is also
448 consistent with previous results using tasks including loss-related options that clearly showed that
449 option valence was taken into account in transfer learning performance (6,8). In addition to that
450 other studies illustrate that in similar paradigms, other behavioral measures, such as reaction times
451 and confidence, are strongly affected by the learning context valence (valence is a construct that is
452 absent in direct policy learning methods)(13,37). Finally, consistent with our intermediate view,
453 other imaging studies found value related representations more consistent with a partial
454 normalization process (38,39).

455

456 **Conclusions**

457 To conclude, we demonstrated that in humans, reinforcement learning values are learnt in a
458 context-dependent manner that is compatible with range adaptation (instantiated as a range
459 normalization process) (40). Specifically, we tested the possibility that these results from the way
460 outcome information is automatically processed to achieve adaptive coding (41), by showing that
461 the lower outcome uncertainty, the fuller range adaptation. This leads to a paradoxical result:
462 reducing task difficulty can, in some occasions, decrease choice optimality. This surprising result
463 can be understood with a perceptual analogy. Going into a dark room forces us to adapt our retinal
464 response to dark, so that when we go back into a light condition we do not see very well. The longer
465 we are exposed to dim light, the stronger the effect when we go back to normal.

466 Our findings fit in the debate aimed at deciding whether the computational processes leading to
467 suboptimal decisions have to be considered flaws or feature of human cognition (42,43). Range
468 adapting reinforcement learning is clearly adaptive in the learning phase. We could hypothesize
469 that the situations in which the process is adaptive are more frequent in real life. In other terms the
470 performance of the system has to be evaluated as a function of the tasks it has been selected to
471 solve. It is true that we may be hit by a bus when we exit a dark room because we do not see well,
472 but on average, the benefit of a sharper perception in a dark room is big enough to compensate for
473 the (rare) event of a bus waiting for us outside the dark room. Ultimately, whether context-

474 dependent reinforcement learning should be considered a flaw or a desirable feature of human
475 cognition should be determined comparing the real life frequency of the situations where it is
476 adaptive (as in the learning phase) to that where it is maladaptive (as in the transfer phase). While
477 our study does not settle this issue, our findings do demonstrate that this process induce, at least in
478 some circumstances, economically suboptimal choices.

479 **Materials and Methods**

480

481 • **Participants**

482

483 For the laboratory experiment, we recruited 40 participants (28 females, aged 24.28 ± 3.05 years
484 old) via Internet advertising in a local mailing-list dedicated to cognitive science-related activities.
485 For the online experiments, we recruited 8x100 participants (414 females, aged 30.06 ± 10.10 years)
486 from the Prolific platform (www.prolific.co). We based the online sample size on a power analysis
487 that was based on the behavioral results of the lab experiment. In the $\Delta EV = 1.75$ context, lab
488 participants reached a difference between choice rate and chance (0.5) of 0.11 ± 0.30
489 (mean \pm s.e.m.). To obtain the same with a power of 0.95, the Matlab function ‘samsizewr.m’
490 indicated a value of 99 participants that we rounded to 100. The research was carried out following
491 the principles and guidelines for experiments including human participants provided in the
492 declaration of Helsinki (1964, revised in 2013). The Inserm Ethical Review Committee /
493 IRB00003888 approved the study on November 13th, 2018 and participants were provided written
494 informed consent prior to their inclusion. To sustain motivation throughout the experiment,
495 participants were given a bonus depending on the number of points won in the experiment (average
496 money won in pounds: 4.14 ± 0.72 , average performance against chance: $M = 0.65 \pm 0.13$,
497 $t(799) = 33.91$, $p < 0.0001$). A laboratory-based experiment was originally performed ($N = 40$) to
498 ascertain that online testing would not significantly affecting the main conclusions. The results are
499 presented in the **Supplementary Materials**.

500

501 • **Behavioral tasks**

502

503 Participants performed an online version of a probabilistic instrumental learning task adapted from
504 previous studies (6). After checking the consent form, participants received written instructions
505 explaining how the task worked and that their final payoff would be affected by their choices in
506 the task. During the instructions the possible outcomes in points (0pt, 1pt and 10pt) were explicitly
507 showed as well as their conversion rate (1pt = 0.005£). The instructions were followed by a short
508 training session of 12 trials aiming at familiarizing the participants with the response modalities.
509 Participants could repeat the training session up to two times and then started the actual experiment.

510

511 In our task, options were materialized by abstract stimuli (cues) taken from randomly generated
512 identicons, colored such that the subjective hue and saturation were very similar according to the
513 HSL_{UV} color scheme (www.hsluv.org). On each trial, two cues were presented on both sides of the
514 screen. The side in which a given cue was presented was pseudo-randomized, such that a given cue
515 was presented an equal number of times on the left and the right. Participants were required to
516 select between the two cues by clicking on one cue. The choice window was self-paced. A brief
517 delay after the choice was recorded (500 ms), the outcome was displayed for 1000 ms. There was
518 no fixation screen between trials. The average reaction time was 1.36 ± 0.04 seconds (median: 1.16),
519 the average experiment completion time was 325.24 ± 8.39 seconds (median: 277.30).

520

521 As in previous studies, the full task consisted in one *learning* phase followed by a *transfer* phase
522 (6–8,44). During the learning phase, cues appeared in four fixed pairs. Each pair was presented 30
523 times, leading to a total of 120 trials. Within each pair, the two cues were associated to a zero and
524 a non-zero outcome with reciprocal probabilities (0.75/0.25 and 0.25/0.75). At the end of the trial,
525 the cues disappeared and the selected one was replaced by the outcome (“10”, “1”, or “0”) (**Figure**

526 **1A)**. In experiments E3, E4, E7 and E8, the outcome corresponding to the forgone option
527 (sometimes referred to as the *counterfactual* outcome) was also displayed (**Figure 1C**). Once they
528 had completed the learning phase, participants were displayed with the total points earned and their
529 monetary equivalent.

530

531 During the transfer phase after the learning phase, the pairs of cues were rearranged into four new
532 pairs. The probability of obtaining a specific outcome remained the same for each cue (**Figure 1B**).
533 Each new pair was presented 30 times, leading to a total of 120 trials. Before the beginning of the
534 transfer phase, participants were explained that they would be presented with the same cues, only
535 that the pairs would not have been necessarily displayed together before. In order to prevent explicit
536 memorizing strategies, participants were not informed that they would have to perform a transfer
537 phase until the end of the learning phase. After making a choice, the cues disappeared. In
538 experiments E1, E3, E5 and E7, participants were not informed of the outcome of the choice on a
539 trial-by-trial basis and the next trial began after 500ms. This was specified in the instruction phase.
540 In experiments E2, E4, E6 and E8, participants were informed about the result of their choices in a
541 trial-by-trial basis and the outcome was presented for 1000ms. In all experiments they were
542 informed about the total points earned at the end of the transfer phase. In addition to the presence
543 / absence of feedback, experiments differed in two other factors. Feedback information could be
544 either partial (experiments E1, E2, E5, E6) or complete, (experiments E3, E4, E7, E8; meaning the
545 outcome of the forgone option was also showed). When the transfer phase included feedback, the
546 information factor was the same as in the learning phase. Trial structure was also manipulated, such
547 that in some experiments (E5, E6, E7, E8), all trials of a given choice context were clustered
548 ('blocked'), and in the remaining experiments (E1, E2, E3, E4) they were interleaved, in both the
549 learning phase and the transfer phase (**Figure 1C**).

550

551

552 **Re-analysis of a previous experiment**

553 In the present paper we also include new analyses of previously published experiments (6). The
554 general constructional principle of the previous experiments is similar to that used in the present
555 experiments, as they involved a learning phase and a transfer phase. However, the previous design
556 different from the present one in several important respects. First, in addition to an outcome
557 magnitude manipulation ('10c vs '1€', similar to our learning phase), this study also manipulated
558 the valence of the outcomes (gain vs loss), generating to a 2x2 factorial design. Second, the
559 organization of the transfer phase was quite different. Indeed, each option was compared with all
560 other possible options. The mean dependent variable extracted from the transfer phase is therefore
561 not the correct response rate, but simply choice rate per option (which is proportional to its
562 subjective value). The data were pooled across two experiments featuring partial (N=20) and partial
563 and complete feedback trials (N=40). In both experiments the choice contexts were interleaved.
564 Other differences include the fact that these previous experiments were laboratory-based and
565 featured a slightly different number of trials, different stimuli and timing (see the original
566 publication for more details).

567

568

- **Analyses**

569 **Behavioral analyses.**

570 The main dependent variable was the *correct* choice rate, i.e., choices directed toward the option
571 with the highest expected value. Statistical effects were assessed using multiple-way repeated
572 measures ANOVAs with choice context (labeled in the manuscript by their difference in expected
573 values: ΔEV) as within-subject factor, and feedback information, feedback in the transfer phase
574 and task structure as between-subjects factors. Post-hoc tests were performed using one-sample
575 and two-sample t-tests for respectively within- and between-experiment comparisons. To assess
576 overall performance, additional one sample t-tests were performed against chance level (0.5). We
577 report the *t*-statistic, *p*-value, and Cohen's *d* to estimate effect size (two-sample t-test only). Given
578 the large sample size (N=800), central limit theorem allows us to assume normal distribution of
579 our overall performance data and to apply properties of normal distribution in our statistical
580 analyses, as well as sphericity hypotheses. Concerning ANOVA analyses, we report the
581 uncorrected statistical, as well as Huynh–Feldt correction for repeated measures ANOVA when
582 applicable (45), *F*-statistic, *p*-value, partial eta-squared η_p^2 and generalized eta-squared η^2 (when
583 Huynh-Feldt correction is applied) to estimate effect size. All statistical analyses were performed
584 using Matlab (www.mathworks.com) and R (www.r-project.org). For visual purposes, learning
585 curves were smoothed using a moving average filter (span of 5 in Matlab's *smooth* function).

586

587 • Models

588 We analyzed our data with variation of simple associative learning models (46,47). The goal of all
589 models is to estimate in each choice context (or *state*) the expected reward (*R*) of each option and
590 pick the one that maximizes this expected reward *R*.

591 At trial *t*, option values of the current context *s* are updated with the delta rule:
592

$$593 Q_{t+1}(s, c) = Q_t(s, c) + \alpha_c \delta_{c,t} \quad (5)$$

$$594 Q_{t+1}(s, u) = Q_t(s, u) + \alpha_u \delta_{u,t} \quad (6)$$

595

596 where α_c is the learning rate for the chosen (*c*) option and α_u the learning rate for the unchosen (*u*)
597 option, i.e., the counterfactual learning rate. δ_c and δ_u are prediction error terms calculated as
598 follows:

$$599 \delta_{c,t} = R_{c,t} - Q_t(s, c) \quad (7)$$

$$600 \delta_{u,t} = R_{u,t} - Q_t(s, u) \quad (8)$$

601

602 δ_c is calculated in both partial and complete feedback contexts and δ_u is calculated in the
603 experiments with complete feedback only.

604

605 We modelled participants' choice behavior using a softmax decision rule representing the
606 probability for a participant to choose one option *a* over the other option *b*:

607

$$608 P_t(s, a) = \frac{1}{1 + e^{(Q_t(s,b) - Q_t(s,a) * \beta)}} \quad (9)$$

609

610 where β is the inverse temperature parameter. High temperatures ($\beta \rightarrow 0$) cause the action to be all
611 (nearly) equiprobable. Low temperatures ($\beta \rightarrow +\infty$) cause a greater difference in selection
612 probability for actions that differ in their value estimates (46).
613

614

615 We compared three alternative computational models: the ABSOLUTE model, which encodes
616 outcomes in an absolute scale independently of the choice context in which they are presented, the
617 RANGE model which tracks the value of the maximum reward in each context and normalizes the
618 actual reward accordingly, rescaling rewards between 0 and 1, and the HABIT model, which
619 integrates action weights into the decision process.

620

621 ABSOLUTE model

622 The outcomes are encoded as the participants see them. A positive outcome is encoded as its actual
623 positive value (in points): $R_{ABS,t} \in \{10, 1, 0\}$.

624

625 RANGE model

626 The outcomes (both chosen and unchosen) are encoded on a context-dependent relative scale. On
627 each trial the relative reward $R_{RAN,t}$ is calculated as follows:

$$628 \quad R_{RAN,t} = \frac{R_{ABS,t} - R_{MIN,t}(s)}{R_{MAX,t}(s) - R_{MIN,t}(s) + 1} \quad (2)$$

629

630 As R_{MIN} is initialized at zero and never changes, in our task this model can be reduced as:

631

$$632 \quad R_{RAN,t} = \frac{R_{ABS,t}}{R_{MAX,t}(s) + 1} \quad (10)$$

633

634 where s is the decision context (i.e., a combination of options) and R_{MAX} is a context-dependent
635 variable, initialized to 0 and updated at each trial t if the outcome is greater than its current value:

636

$$637 \quad R_{MAX,t+1}(s) = R_{MAX,t}(s) + \alpha_R (R_{ABS,t} - R_{MAX,t}(s)) \quad \text{if } R_{ABS,t} > R_{MAX,t}(s) \quad (11)$$

638

639 Accordingly, outcomes are progressively normalized so that eventually $R_{RAN,t} \in [0,1]$. The
640 chosen and unchosen option values and prediction errors are updated with the same rules as in the
641 ABSOLUTE model. Note that the ABSOLUTE model is nested within the RANGE model ($\alpha_R =$
642 0).

643

644 HABIT model

645 The outcomes are encoded on an absolute scale, but decisions integrate a habitual component
646 (16,19). To do so, in addition to the Q-values, a habitual (or choice trace) component H is tracked
647 and updated (with a dedicated learning rate parameter) that takes into account the selected action
648 (1 for chosen option, 0 for the unchosen option). The choice is performed with a softmax rule based
649 on decision weights D that integrate Q-values and decision weights H :

$$650 \quad D_t(s, c) = (1 - \omega) * Q_t(s, c) + \omega * H_t(s, c) \quad (3)$$

651 where at each trial t , state s and chose option c , D is the arbiter, Q is the goal directed component
652 (Q-values matrix), H is the habitual component. The weight ω is fitted as an additional parameter
653 and governs the relative weights of values and habits (for $\omega=0$ the model reduces to the
654 ABSOLUTE model).

655

656 UTILITY model

657 The outcomes are encoded as an exponentiation of the absolute reward, leading to a curvature of
658 the value function (29):

659
$$R_{\text{UTI},t} = (R_{\text{ABS},t})^\nu \quad (4)$$

660 where the exponent ν is the utility parameter, with $0 < \nu < 1$ (for $\nu = 1$ the model reduces to the
661 ABSOLUTE model).

662 **References and Notes**

- 663 1. Louie K, Glimcher PW. Efficient coding and the neural representation of value. *Ann N Y Acad Sci.* 2012
664 Mar;1251:13–32.
- 665 2. Vlaev I, Chater N, Stewart N, Brown GDA. Does the brain calculate value? *Trends Cogn Sci.* 2011
666 Nov;15(11):546–54.
- 667 3. Cox KM, Kable JW. BOLD Subjective Value Signals Exhibit Robust Range Adaptation. *J Neurosci.* 2014 Dec
668 3;34(49):16533–43.
- 669 4. Nieuwenhuis S, Heslenfeld DJ, Alting von Geusau NJ, Mars RB, Holroyd CB, Yeung N. Activity in human
670 reward-sensitive brain areas is strongly context dependent. *NeuroImage.* 2005 May 1;25(4):1302–9.
- 671 5. Elliott R, Agnew Z, Deakin JFW. Medial orbitofrontal cortex codes relative rather than absolute value of
672 financial rewards in humans. *Eur J Neurosci.* 2008 May;27(9):2213–8.
- 673 6. Bavard S, Lebreton M, Khamassi M, Coricelli G, Palminteri S. Reference-point centering and range-
674 adaptation enhance human reinforcement learning at the cost of irrational preferences. *Nat Commun.* 2018
675 29;9(1):4503.
- 676 7. Klein TA, Ullsperger M, Jocham G. Learning relative values in the striatum induces violations of normative
677 decision making. *Nat Commun.* 2017 Jun 20;8(1):16033.
- 678 8. Palminteri S, Khamassi M, Joffily M, Coricelli G. Contextual modulation of value signals in reward and
679 punishment learning. *Nat Commun.* 2015 Aug 25;6:8096.
- 680 9. Freidin E, Kacelnik A. Rational Choice, Context Dependence, and the Value of Information in European
681 Starlings (*Sturnus vulgaris*). *Science.* 2011 Nov 18;334(6058):1000–2.
- 682 10. Pompilio L, Kacelnik A. Context-dependent utility overrides absolute memory as a determinant of choice.
683 *Proc Natl Acad Sci.* 2010 Jan 5;107(1):508–12.
- 684 11. Rustichini A, Conen KE, Cai X, Padoa-Schioppa C. Optimal coding and neuronal adaptation in economic
685 decisions. *Nat Commun.* 2017 Oct 31;8(1):1208.
- 686 12. Webb R, Glimcher PW, Louie K. The Normalization of Consumer Valuations: Context-Dependent
687 Preferences From Neurobiological Constraints. *Manag Sci* [Internet]. 2020 May 27 [cited 2020 Jul 27];
688 Available from: <https://pubsonline.informs.org/doi/abs/10.1287/mnsc.2019.3536>
- 689 13. Fontanesi L, Palminteri S, Lebreton M. Decomposing the effects of context valence and feedback information
690 on speed and accuracy during reinforcement learning: a meta-analytical approach using diffusion decision
691 modeling. *Cogn Affect Behav Neurosci.* 2019 Jun 1;19(3):490–502.
- 692 14. Collins AGE, Frank MJ. How much of reinforcement learning is working memory, not reinforcement
693 learning? A behavioral, computational, and neurogenetic analysis. *Eur J Neurosci.* 2012 Apr;35(7):1024–35.
- 694 15. Rabin M. Diminishing Marginal Utility of Wealth Cannot Explain Risk Aversion. 2000 Jun 9 [cited 2020 Jul
695 27]; Available from: <https://escholarship.org/uc/item/61d7b4pg>
- 696 16. Miller KJ, Shenhav A, Ludvig EA. Habits without values. *Psychol Rev.* 2019;126(2):292–311.
- 697 17. Landry P, Webb R. Pairwise Normalization: A Neuroeconomic Theory of Multi-Attribute Choice [Internet].
698 Rochester, NY: Social Science Research Network; 2019 Jan [cited 2020 Nov 30]. Report No.: ID 2963863.
699 Available from: <https://papers.ssrn.com/abstract=2963863>

- 700 18. Palminteri S, Wyart V, Koehlin E. The Importance of Falsification in Computational Cognitive Modeling. Trends Cogn Sci. 2017 Jun;21(6):425–33.
701
- 702 19. Katahira K. The statistical structures of reinforcement learning with asymmetric value updates. J Math
703 Psychol. 2018 Dec 1;87:31–45.
- 704 20. Louie K, Glimcher PW, Webb R. Adaptive neural coding: from biological to behavioral decision-making. Curr
705 Opin Behav Sci. 2015 Oct 1;5:91–9.
- 706 21. Dumbalska T, Li V, Tsetsos K, Summerfield C. A map of decoy influence in human multialternative choice.
707 Proc Natl Acad Sci. 2020 Oct 6;117(40):25169–78.
- 708 22. Daviet R, Webb R. A Double Decoy Experiment to Distinguish Theories of Dominance Effects [Internet].
709 Rochester, NY: Social Science Research Network; 2019 Mar [cited 2020 Nov 30]. Report No.: ID 3374514.
710 Available from: <https://papers.ssrn.com/abstract=3374514>
- 711 23. Gluth S, Kern N, Kortmann M, Vitali CL. Value-based attention but not divisive normalization influences
712 decisions with multiple alternatives. Nat Hum Behav. 2020 Jun;4(6):634–45.
- 713 24. Goodwin PB. Habit and Hysteresis in Mode Choice. Urban Stud. 1977;14(1):95–8.
- 714 25. Dickinson A, Weiskrantz L. Actions and habits: the development of behavioural autonomy. Philos Trans R
715 Soc Lond B Biol Sci. 1985 Feb 13;308(1135):67–78.
- 716 26. Lally P, Jaarsveld CHM van, Potts HWW, Wardle J. How are habits formed: Modelling habit formation in the
717 real world. Eur J Soc Psychol. 2010;40(6):998–1009.
- 718 27. Thrailkill EA, Trask S, Vidal P, Alcalá JA, Bouton ME. Stimulus Control of Actions and Habits: A Role for
719 Reinforcer Predictability and Attention in the Development of Habitual Behavior. J Exp Psychol Anim Learn
720 Cogn. 2018 Oct;44(4):370–84.
- 721 28. Akaishi R, Umeda K, Nagase A, Sakai K. Autonomous mechanism of internal choice estimate underlies
722 decision inertia. Neuron. 2014 Jan 8;81(1):195–206.
- 723 29. Neumann J von, Morgenstern O. Theory of Games and Economic Behavior. Princeton University Press; 1953.
724 774 p.
- 725 30. Bernoulli D. Exposition of a New Theory on the Measurement of Risk. Econometrica. 1954;22(1):23–36.
- 726 31. Markowitz H. The Utility of Wealth. J Polit Econ. 1952 Apr 1;60(2):151–8.
- 727 32. Kahneman D, Tversky A. Subjective probability: A judgment of representativeness. Cognit Psychol. 1972 Jul
728 1;3(3):430–54.
- 729 33. Loomes G, Sugden R. Regret Theory: An Alternative Theory of Rational Choice Under Uncertainty. Econ J.
730 1982;92(368):805–24.
- 731 34. Dayan P, Abbott LF. Theoretical Neuroscience: Computational and Mathematical Modeling of Neural
732 Systems. Massachusetts Institute of Technology Press; 2001. 460 p.
- 733 35. Li J, Daw ND. Signals in human striatum are appropriate for policy update rather than value prediction. J
734 Neurosci Off J Soc Neurosci. 2011 Apr 6;31(14):5504–11.
- 735 36. Palminteri S, Pessiglione M. Chapter 23 - Opponent Brain Systems for Reward and Punishment Learning:
736 Causal Evidence From Drug and Lesion Studies in Humans. In: Dreher J-C, Tremblay L, editors. Decision

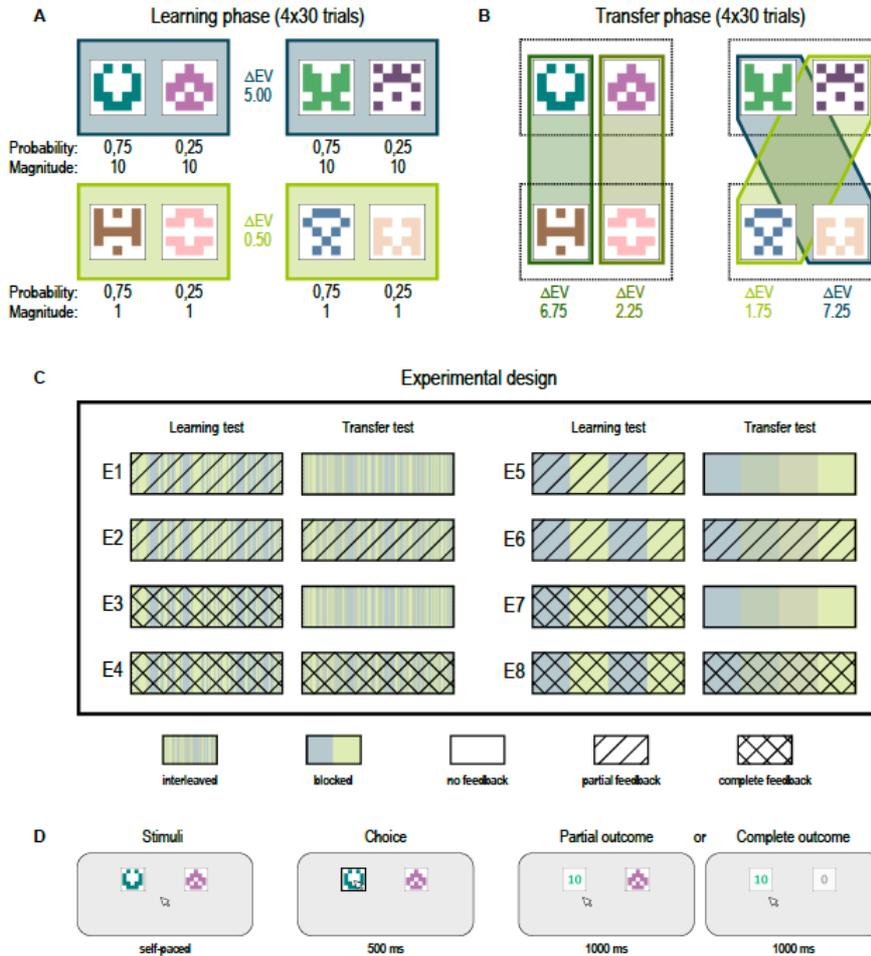
- 737 Neuroscience [Internet]. San Diego: Academic Press; 2017 [cited 2020 Nov 30]. p. 291–303. Available from:
738 <http://www.sciencedirect.com/science/article/pii/B9780128053089000233>
- 739 37. Lebreton M, Bacily K, Palminteri S, Engelmann JB. Contextual influence on confidence judgments in human
740 reinforcement learning. *PLOS Comput Biol*. 2019 avr;15(4):e1006973.
- 741 38. Burke CJ, Baddeley M, Tobler PN, Schultz W. Partial Adaptation of Obtained and Observed Value Signals
742 Preserves Information about Gains and Losses. *J Neurosci*. 2016 Sep 28;36(39):10016–25.
- 743 39. Pischedda D, Palminteri S, Coricelli G. The Effect of Counterfactual Information on Outcome Value Coding
744 in Medial Prefrontal and Cingulate Cortex: From an Absolute to a Relative Neural Code. *J Neurosci*. 2020 Apr
745 15;40(16):3268–77.
- 746 40. Conen KE, Padoa-Schioppa C. Partial Adaptation to the Value Range in the Macaque Orbitofrontal Cortex. *J*
747 *Neurosci*. 2019 May 1;39(18):3498–513.
- 748 41. Padoa-Schioppa C, Rustichini A. Rational Attention and Adaptive Coding: A Puzzle and a Solution. *Am Econ*
749 *Rev*. 2014 May;104(5):507–13.
- 750 42. Gigerenzer G. The Bias Bias in Behavioral Economics. *Rev Behav Econ*. 2018 Dec 30;5(3–4):303–36.
- 751 43. Haselton MG, Nettle D, Andrews PW. The Evolution of Cognitive Bias. In: *The Handbook of Evolutionary*
752 *Psychology* [Internet]. John Wiley & Sons, Ltd; 2015 [cited 2020 Jul 27]. p. 724–46. Available from:
753 <https://onlinelibrary.wiley.com/doi/abs/10.1002/9780470939376.ch25>
- 754 44. Frank MJ, Seeberger LC, O'reilly RC. By carrot or by stick: cognitive reinforcement learning in parkinsonism.
755 *Science*. 2004 Dec 10;306(5703):1940–3.
- 756 45. Girden ER. ANOVA: Repeated Measures. SAGE; 1992. 88 p.
- 757 46. Sutton RS, Barto AG. Reinforcement Learning - An Introduction [Internet]. Mit Press; 1998 [cited 2017 Jun
758 8]. Available from: <http://gen.lib.rus.ec/book/index.php?md5=C5BC41EB7F60C05D37473B4A9AC77A2A>
- 759 47. Rescorla RA, Wagner AR. A theory of Pavlovian conditioning: Variations in the effectiveness of
760 reinforcement and nonreinforcement. *Class Cond II Curr Res Theory*. 1972;2:64–99.
- 761 48. Hergueux J, Jacquemet N. Social preferences in the online laboratory: a randomized experiment. *Exp Econ*.
762 2015 Jun 1;18(2):251–83.
- 763 49. Kahneman D. Maps of Bounded Rationality: Psychology for Behavioral Economics. *Am Econ Rev*. 2003
764 Dec;93(5):1449–75.
- 765 50. Shavit T, Sonsino D, Benzion U. A comparative study of lotteries-evaluation in class and on the Web. *J Econ*
766 *Psychol*. 2001 Aug 1;22(4):483–91.
- 767 51. Schoeffler M, Stöter F-R, Bayerlein H, Edler B, Herre J. An Experiment about Estimating the Number of
768 Instruments in Polyphonic Music: A Comparison Between Internet and Laboratory Results. In: *ISMIR*. 2013.
- 769 52. Reinecke K, Gajos KZ. LabintheWild: Conducting Large-Scale Online Experiments With Uncompensated
770 Samples. In: *Proceedings of the 18th ACM Conference on Computer Supported Cooperative Work & Social*
771 *Computing* [Internet]. Vancouver, BC, Canada: Association for Computing Machinery; 2015 [cited 2020 Jul
772 27]. p. 1364–1378. (CSCW '15). Available from: <https://doi.org/10.1145/2675133.2675246>

773

774 Acknowledgments

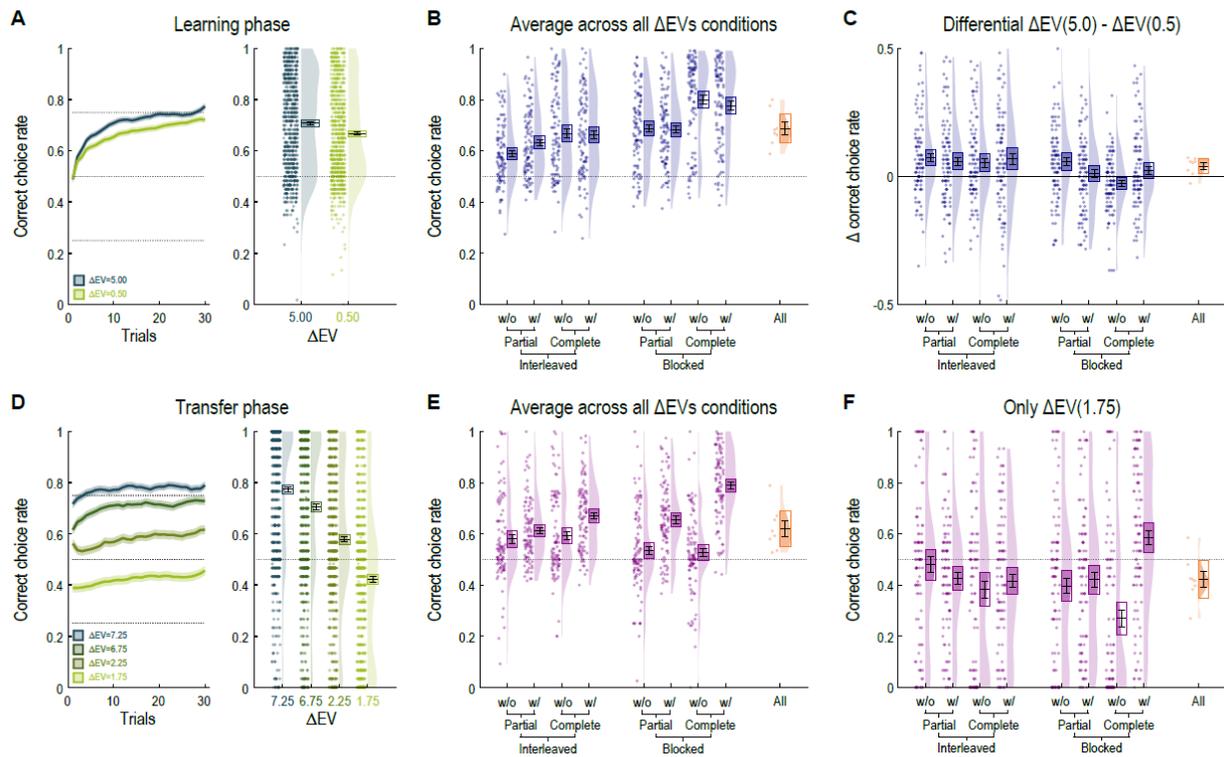
775 **Funding:** S.P. is supported by an ATIP-Avenir grant (R16069JS), the Programme Emergence(s)
776 de la Ville de Paris, the Fondation Fyssen and the Fondation Schlumberger pour l'Education et la
777 Recherche. S.B. is supported by MILDECA (Mission Interministérielle de Lutte contre les Drogues
778 et les Conduites Addictives) and the EHESS (Ecole des Hautes Etudes en Sciences Sociales). The
779 funding agencies did not influence the content of the manuscript.

780 **Figures and Tables**



781
782
783
784
785
786
787
788
789
790
791
792

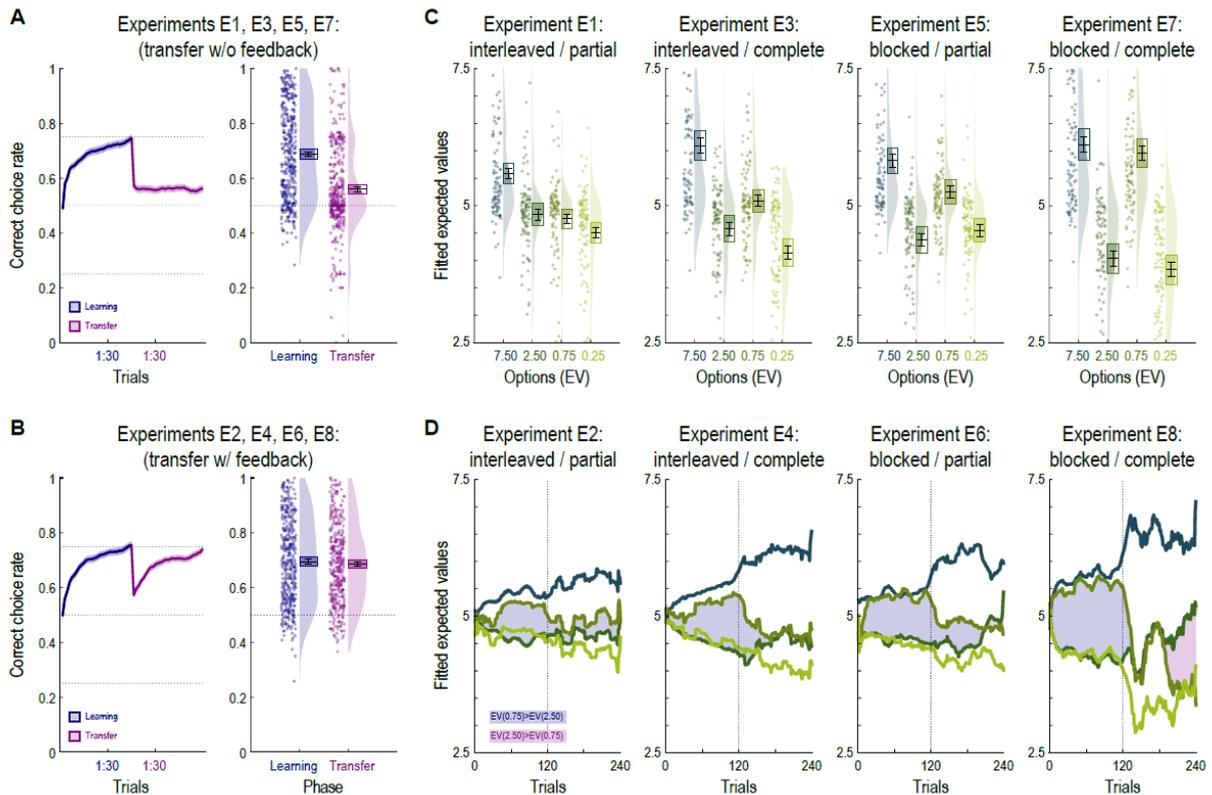
Fig. 1. Experimental design. (A) Choice contexts in the learning phase. During the learning phase, participants were presented with four choice contexts, including high magnitude ($\Delta EV = 5.0$ contexts) and low magnitude ($\Delta EV = 0.5$ contexts). (B) Choice contexts in the transfer phase. The four options were re-arranged into four new choice contexts, each involving both the 1pt and the 10pt outcome. (C). Experimental design. The eight experiments varied in the temporal arrangement of choice contexts (interleaved or block) and the quantity of feedback in the learning phase (partial or complete) and the transfer phase (without or with). (D) Successive screens of a typical trials (complete feedback; durations are given in milliseconds).



793

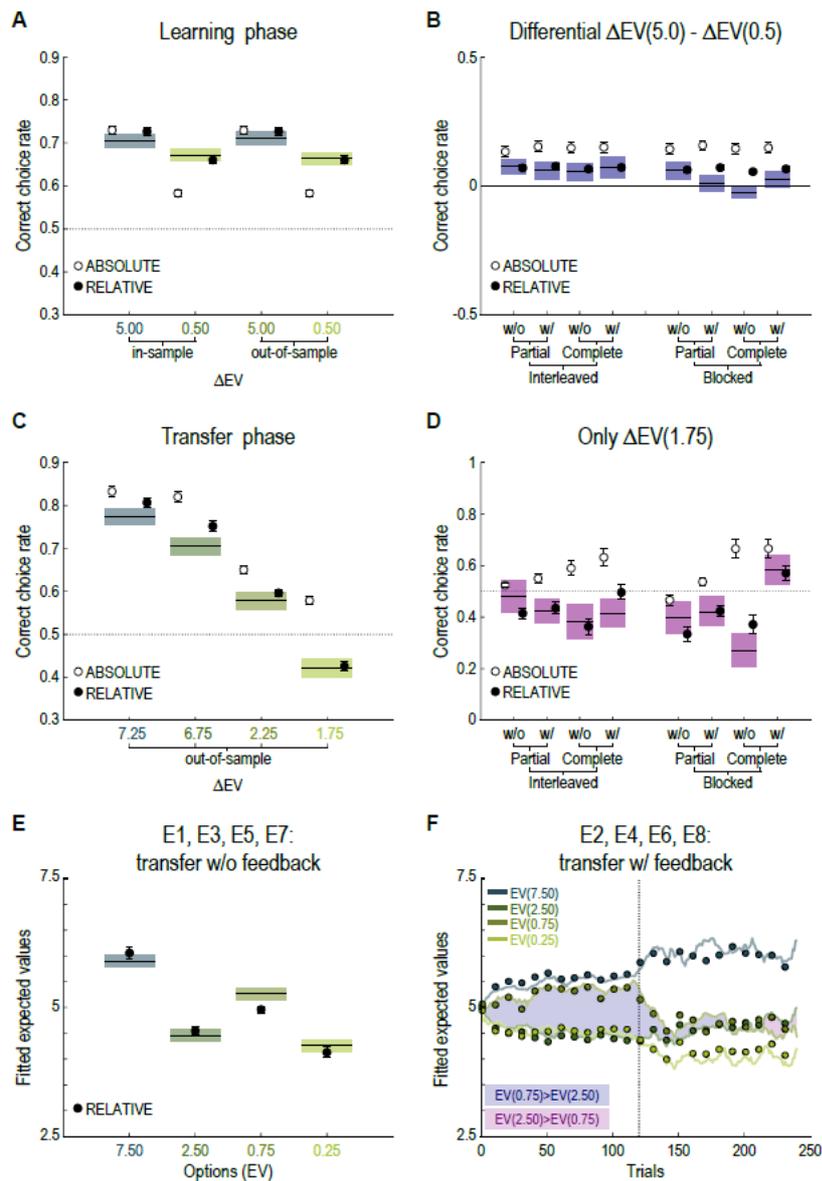
794 **Figure 2. Behavioral results.** (A) Correct choice rate in the learning phase as a function of the
 795 choice context ($\Delta EV=5.0$ or $\Delta EV=0.5$). Leftmost panel: learning curves; rightmost panel: average
 796 across all trials. (B) Average correct response rate in the learning phase per experiment (in blue:
 797 $N=800$ participants) and meta-analytical (in orange: $N=8$ experiments). (C) Difference in correct
 798 choice rate between the $\Delta EV=5.0$ and the $\Delta EV=0.5$ contexts per experiment (in blue: $N=800$
 799 participants) and meta-analytical (in orange: $N=8$ experiments). (D) Correct choice rate in the
 800 transfer phase as a function of the choice context ($\Delta EV=7.25$, $\Delta EV=6.75$, $\Delta EV=2.25$ or
 801 $\Delta EV=1.75$). Leftmost panel: learning curves; rightmost panel: average across all trials. (E) Average
 802 correct response rate in the transfer phase per experiment (in pink: $N=800$ participants) and meta-
 803 analytical (in orange: $N=8$ experiments). (F) Correct choice rate for the $\Delta EV=1.75$ context only (in
 804 pink: $N=800$ participants) and meta-analytical (in orange: $N=8$ experiments). In all panels: points
 805 indicate individual average, areas indicate probability density function, boxes indicate 95%
 806 confidence interval and errors bars indicate s.e.m.

807



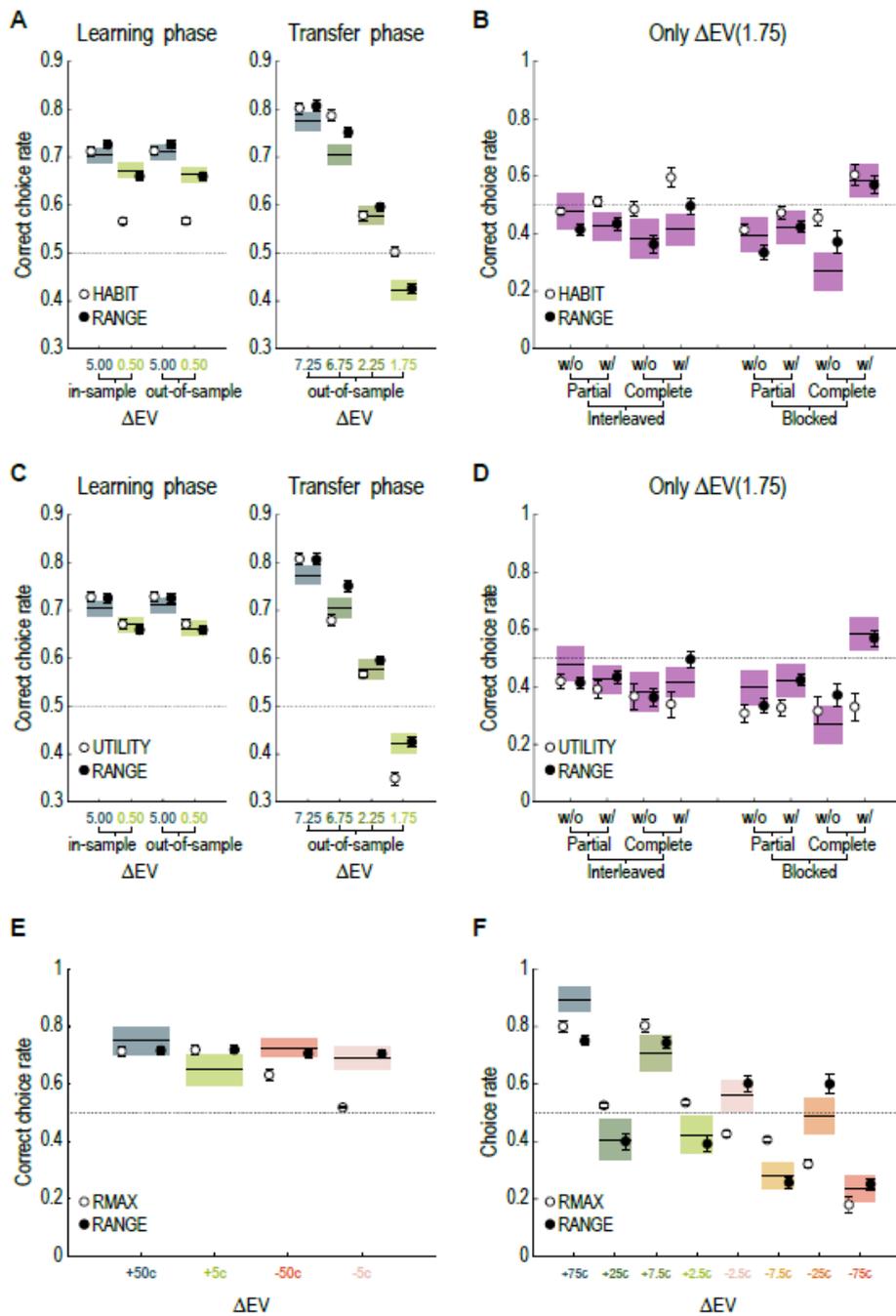
808
809
810

811 **Figure 3. Learning versus transfer comparison and inferred option values. (A-B)** Average
812 response rate in the learning (blue) and transfer (pink) phase for experiments without (A) and with
813 (B) trial-by-trial transfer feedback. Leftmost panel: learning curves; rightmost panel: average
814 across all trials. (C) Average inferred option values for the experiments without trial-by-trial
815 transfer feedback. (D) Trial-by-trial inferred option values for the experiments with trial-by-trial
816 transfer feedback. In all panels: points indicate individual average, areas indicate probability
817 density function, boxes indicate 95% confidence interval and errors bars indicate s.e.m.
818



819
820

821 **Figure 4. Model comparison.** Model simulations of ABSOLUTE (white) and RANGE (black)
 822 models over the behavioral data (mean and 95% confidence interval) in each context. **(A)**
 823 Simulated data in the learning phase were obtained with the parameters fitted in half the contexts
 824 ($\Delta EV=5.0$ and $\Delta EV=0.5$) of the learning phase (in-sample and out-of-sample predictions). **(B)**
 825 Data and simulations of the differential between high magnitude ($\Delta EV=5.0$) and low magnitude
 826 ($\Delta EV=0.5$) contexts. **(C)** Simulated data in the transfer phase were obtained with the parameters
 827 fitted in all the contexts of the learning phase (out-of-sample predictions). **(D)** Data and simulations
 828 in the context $\Delta EV=1.75$ only. **(E)** Average inferred option values for the behavioral data and
 829 simulated data (black dots: RANGE model only) for the experiments without trial-by-trial transfer
 830 feedback. **(F)** Trial-by-trial inferred option values for the behavioral data and simulated data
 831 (colored dots: RANGE model only) for the experiments with trial-by-trial transfer feedback, where
 832 curves indicate trial-by-trial fit of each inferred option value, and colored dots indicate RANGE
 833 model simulations.

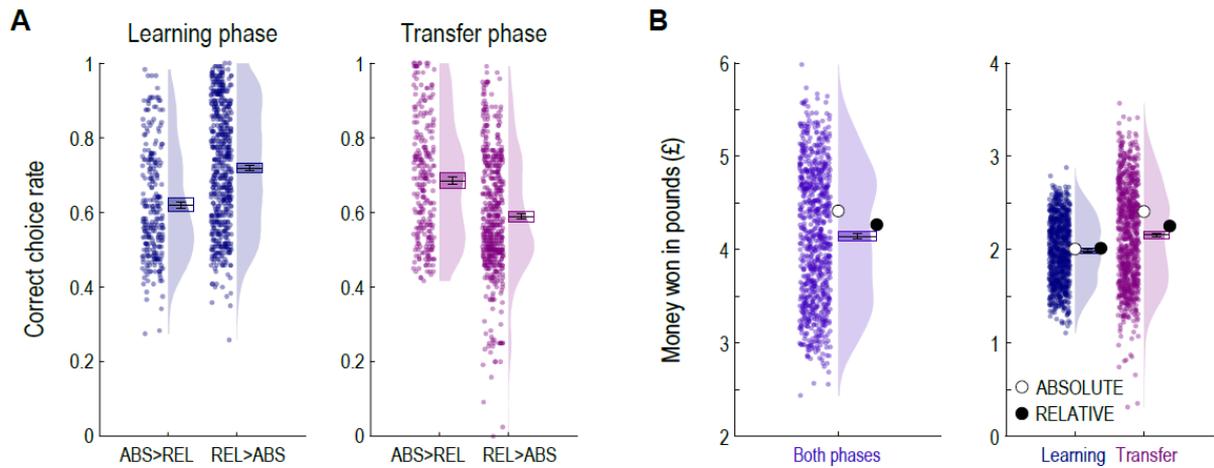


834
835
836
837
838
839
840
841
842
843
844

Figure 5. Ruling out alternative models. Model simulations of HABIL, resp. UTILITY (white) and RANGE (black) models over the behavioral data (mean and 95% confidence interval) in each context. (A, C) Simulated data in the learning phase were obtained with the parameters fitted in half the contexts ($\Delta EV=5.0$ and $\Delta EV=0.5$) of the learning phase (in-sample and out-of-sample predictions). Simulated data in the transfer phase were obtained with the parameters fitted in all the contexts of the learning phase (out-of-sample predictions). (B, D) Data and simulations in the context $\Delta EV=1.75$ only. (E, F) Behavioral data from Bavard, Lebreton et al (2018). Comparing the full RANGE model to its simplified version RMAX in the learning phase (correct choice rate per choice context) and in the transfer test (choice rate per symbol). This study included gain -related

845 contexts (with +50c and +5c as average outcomes) and loss-related contexts (with -50c and -5c as
846 average outcomes) in the learning phase. Choice rates in the transfer phase are presented as a
847 function of decreasing option expected values.

848
849



850
851
852
853
854
855
856
857
858
859
860

Figure 6. The financial cost of relative value learning. (A) Correct choice rate in the learning phase (blue) and the transfer phase (pink) as a function of the difference in log-likelihood between the ABSOLUTE and the RANGE models. ABS>RAN: positive difference, N=255. RAN>ABS: negative difference, N=545. (B) Actual and simulated money won in pounds over the whole task (purple), the learning phase only (blue) and the transfer phase only (pink). Points indicate individual participants, areas indicate probability density function, boxes indicate confidence interval and errors bars indicate s.e.m. Dots indicate model simulations of ABSOLUTE (white) and RANGE (black) models.

D F n	DFd	Learning performance				Transfer performance				Overall performance				
		F-val	Diff	η_p^2		F-val	Diff	η_p^2		F-val	Diff	η_p^2		
L F - L e a r n i n g f e e d b a c k C o m p l e t e > P a r t i a l	1	792	55,57	***	0,079	0,18	22,36	***	0,050	0,01	61,68	***	0,064	0,11
T F - T r a n s f e r f e e d b a c k W i t h > W i t h o u t	1	792	0,04		0,002	0,00	137,18	***	0,12	0,07	58,11	***	0,063	0,10
B E - B l o c k e f f e	1	792	87,22	***	0,099	0,25	1,53		0,013	0,00	46,82	***	0,056	0,08

c t B l o c k > I n t e r l e a v e d P E - P h a s e e f f e c t L e a r n i n g > T r a n s f e r L F x T F L F x B E T F x B E L F x P E T F x P E B E x P E L F x T F														
	1	792	-	-	-	-	-	-	103,07	***	0,067	0,12		
	1	792	2,61		0,01	20,18	***	0,01	3,33			0,01		
	1	792	5,05	*	0,02	1,66		0,00	5,20	*		0,01		
	1	792	2,43		0,01	42,22	***	0,02	9,89	**		0,02		
	1	792	-	-	-	-	-	-	4,97	*		0,01		
	1	792	-	-	-	-	-	-	82,30	***		0,09		
	1	792	-	-	-	-	-	-	42,09	***		0,05		
	1	792	0,55		0,00	5,02	*	0,00	3,65			0,01		

	Experiment 1 N=100		Experiment 2 N=100		Experiment 3 N=100		Experiment 4 N=100		Experiment 5 N=100		Experiment 6 N=100		Experiment 7 N=100		Experiment 8 N=100		Total N=800	
	mean	std	mean	std	mean	std	mean	std	mean	std	mean	std	mean	std	mean	std	mean	std
Age	30,4 8	10,7 0	27,2 3	8,3 0	32,0 1	10,5 1	31,5 7	9,8 0	33,0 4	10,4 8	28,4 6	10,2 0	28,7 3	9,8 9	28,8 4	9,6 0	30,0 6	10,1 0
% correct																		
Learning phase	0,59	0,12	0,63	0,1 3	0,67	0,17	0,66	0,1 6	0,69	0,15	0,68	0,14	0,80	0,1 7	0,78	0,1 6	0,69	0,16
ΔEV=5.0	0,63	0,16	0,66	0,1 7	0,70	0,19	0,70	0,2 0	0,72	0,17	0,69	0,17	0,79	0,1 9	0,79	0,1 8	0,67	0,18
ΔEV=0.5	0,55	0,13	0,60	0,1 4	0,64	0,19	0,63	0,1 9	0,66	0,17	0,68	0,14	0,81	0,1 7	0,76	0,1 8	0,71	0,18
Transfer phase	0,58	0,17	0,61	0,1 2	0,59	0,16	0,67	0,1 3	0,54	0,16	0,66	0,14	0,53	0,1 6	0,79	0,1 4	0,62	0,17
ΔEV=7.25	0,67	0,28	0,76	0,2 2	0,75	0,29	0,85	0,1 9	0,66	0,30	0,84	0,18	0,76	0,3 1	0,93	0,1 4	0,77	0,26
ΔEV=6.75	0,64	0,29	0,68	0,2 6	0,70	0,31	0,81	0,1 9	0,62	0,32	0,76	0,27	0,55	0,3 7	0,89	0,1 6	0,71	0,30
ΔEV=2.25	0,54	0,27	0,58	0,1 9	0,54	0,34	0,61	0,2 8	0,47	0,32	0,60	0,18	0,54	0,3 6	0,76	0,2 2	0,58	0,29
ΔEV=1.75	0,48	0,30	0,43	0,2 3	0,38	0,33	0,42	0,2 7	0,40	0,31	0,42	0,28	0,27	0,3 2	0,59	0,2 9	0,42	0,30

869
870
871
872
873
874
875
876

Table 2. Participants' age and correct choice rate as a function of experiments and task factors.

Model	Learning phase	Transfer phase
ABSOLUTE	-42.74 ± 1.27***	-161.19 ± 11.41***
RANGE	-37.72 ± 0.96	-96.79 ± 4.79
HABIT	-36.68 ± 0.91	-104.62 ± 6.01 ^s
UTILITY	-36.31 ± 0.53***	-104.94 ± 5.24**

877
878
879
880
881
882
883
884
885

Table 3. Quantitative model comparing. Values reported here represent out-of-sample likelihood after two-fold cross validation. Comparison to the RANGE model: *p<0.001; **p<0.01; ^sp<0.08.**

886 **Supplementary Materials**

887

888

Laboratory-based replications and robustness to outliers

889

890 To ascertain that online testing would not significantly affect the main conclusions, a laboratory-
891 based experiment was originally performed. We recruited 40 participants (28 females, aged
892 24.28 ± 3.05) via Internet advertising in mailing list dedicated to cognitive science-related activities.
893 The experimental design was that of experiment E2 of the main text (partial feedback information
894 in both the learning phase and the transfer phase, and trials in an interleaved order; see **Figure 1**).

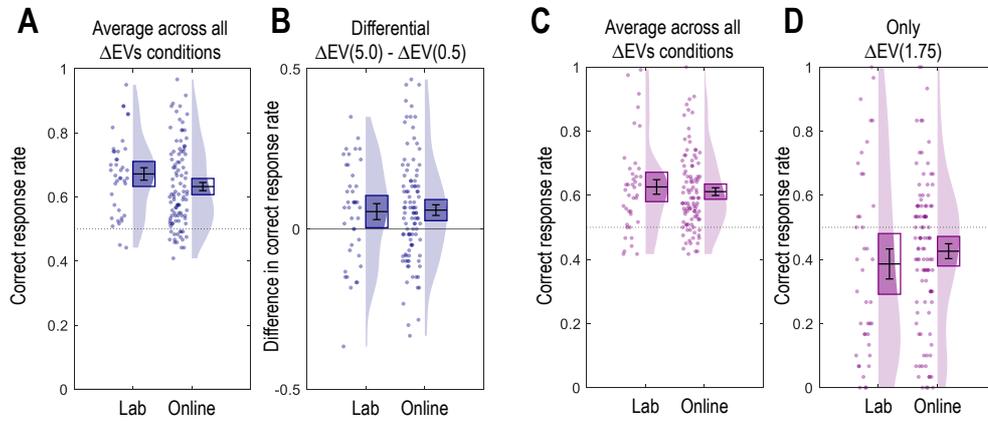
895 In order to characterize learning behavior of participants, we analyzed the correct response rate in
896 both phases, i.e., choices directed toward the most favorable option at each trial. To assess
897 successful learning, we first tested participants' correct response rate against chance level. We
898 found it to be above chance level in both the learning phase ($t(39) = 8.88, p < .0001, d = 1.40,$
899 **Supp. Figure 1A**) and the transfer phase ($t(39) = 5.55, p < .0001, d = 0.88,$ **Supp. Figure 1C**). We
900 found a significant effect of magnitude in the learning phase ($t(39) = 2.18, p = .036, d = 0.34,$ **Supp.**
901 **Figure 1B**), and the correct choice rate in the $\Delta EV = 1.75$ context was significantly below chance
902 level ($t(39) = -2.43, p = .020, d = -0.38,$ **Supp. Figure 1D**). Of note, the effect sizes were really
903 comparable with the ones observed in the corresponding online experiment (learning performance
904 $d = 1.04$ vs 1.40 , transfer performance $d = 0.93$ vs 0.88 , magnitude effect $d = 0.35$ vs 0.34 , value
905 inversion $d = -0.32$ vs -0.38).

906 In addition to checking that the same significant results were present, to formally assess the
907 similarity between online- and laboratory-based experiments, we explicitly compared their scores.
908 Correct choice rate in the learning phase did not significantly differ between laboratory and online
909 datasets ($t(138) = 1.67, p = .10, d = 0.31,$ **Supp. Figure 1A**), neither did the magnitude effect
910 ($t(138) = -0.15, p = .88, d = -0.03,$ **Supp. Figure 1B**). Concerning the transfer phase, overall correct
911 choice rate was not significantly different ($t(138) = 0.62, p = .54, d = 0.12,$ **Supp. Figure 1C**) and
912 the same result was obtained looking specifically at the $\Delta EV = 1.75$ context ($t(138) = -0.84, p = .40,$
913 $d = -0.16,$ **Supp. Figure 1D**). Of note, although the control over the measure of reaction times is
914 arguably limited in online experiments, also this measure did not differ between laboratory- and
915 online-based experiments ($t(138) = -0.50, p = .62, d = -0.09$). This similarity between laboratory-
916 and online-based results supports the usefulness of online-based experiments as a way to target
917 larger, more diversified populations with reduced administrative and financial costs (48). The
918 limitations that can be encountered with online-based experiments - such as lower data quality,
919 faster reaction time, lack of engagement from the participants (49,50) - were not significantly
920 present in our data.

921

922 However, to further check the robustness of our results, we run analyses of the online data
923 excluding subjects presenting unusual task completion time. We approximated participants' total
924 reaction time over the whole task by a normal distribution and removed outliers at a significance
925 level of $p < 0.05$. This led to a removal of 30 participants for the eight online experiments leading
926 to a final sample of 770 participants. We found that the totality of the statistically significant results
927 described in the **Results** section were observable without these reaction time outliers, thus we
928 decided to include all participants in the statistics reported in the **Results** section. In conclusion,
929 our results successfully replicate in the laboratory. Moreover, our results confirm the findings of

930 recent studies comparing both experimental methods and showing that they produce comparable
931 data quality (51,52).
932
933
934



935
936 **Supp. Figure 1. Comparing laboratory and online experiments.** (A) Average correct response
937 rate in the learning phase per experiment. (B) Difference in correct choice rate between the
938 $\Delta EV=5.0$ and the $\Delta EV=0.5$ contexts. (C) Average correct response rate in the transfer phase. (D)
939 Correct choice rate for the $\Delta EV=1.75$ context only.

940
941

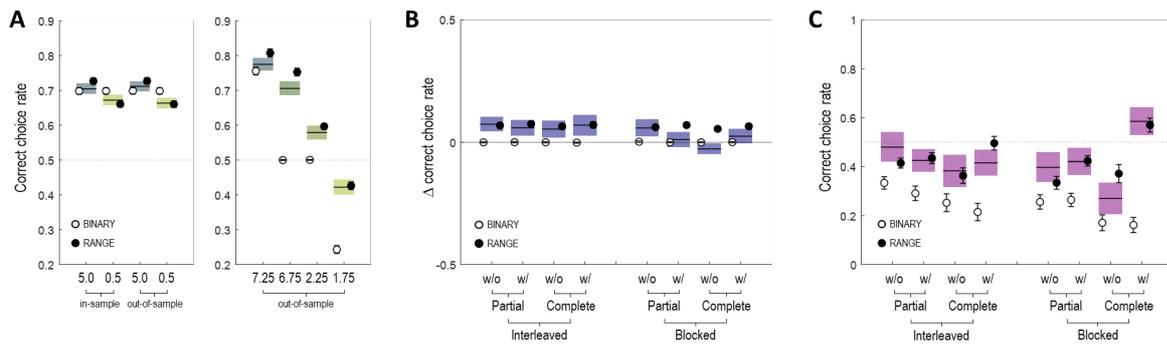
942 **Additional model comparisons**

943 The computational results presented here follow the same fitting and simulation methods presented
 944 in the main text for the main computational models.

945

946 **BINARY model**

947 We analyzed the generative performances of a “full-adaptation” model encoding non-zero
 948 outcomes as ones, regardless of their actual magnitude (10pt, 1pt), that we refer to as the BINARY
 949 model. At least three behavioral features allow us to reject the BINARY model. Of note, the modek
 950 is a special case of the UTILITY model for extremely diminishing marginal utility ($v = 0$; $R_{UTI,t} =$
 951 $(R_{ABS,t})^v$). First, it is not able to capture participants’ behavior in the learning phase by failing to
 952 accurately predict the outcome magnitude difference (**Supp. Figure 2A** and **Supp. Figure 2B**);
 953 second, the model predicts perfect indifference in the $\Delta EV=6.75$ and the $\Delta EV=2.25$ contexts in the
 954 transfer phase, while behavioral results show, respectively, a strong and moderate preference for
 955 the high EV options in these contexts; third, the BINARY model predicts an exaggerated rate of
 956 suboptimal preferences in the $\Delta EV=1.75$ context in the transfer phase (**Supp. Figure 2A** and **Supp.**
 957 **Figure 2C**). This is true in all 8 experiments and even more striking in E8 where the participants
 958 were able to correct their bias.



959 **Supp. Figure 2:** generative performance of the RANGE model (black dots) compared to a full-
 960 adaptation model encoding rewards as 1’s or 0’s (white dots: BINARY model). Black lines
 961 represent the empirical averages. Colored squares indicate the s.e.m. around the empirical averages.
 962
 963
 964

965 **REFERENCE model**

966 We also analyzed the generative performances of a previous context-dependence model (8) that we
 967 call here REFERENCE because of its distinctive feature is to apply reference point dependence to
 968 outcome encoding:

969

970

971

$$Q(s, a) \leftarrow Q(s, a) + \alpha_Q * (R_{ABS} - V(s) - Q(s, a))$$

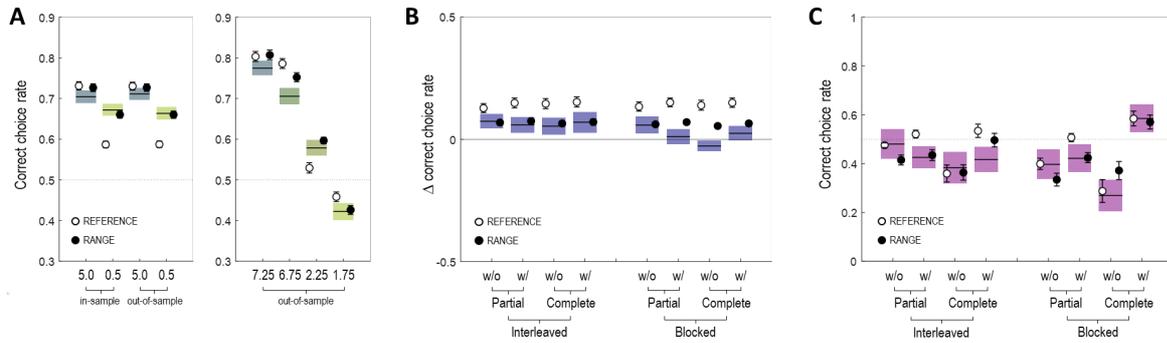
972 Where s is the state (or context, pair of options), $V(s)$ is the state value (or reference point), $Q(s, a)$
 973 is the Q-value (estimated expected value). $V(s)$ is also learnt iteratively, as follows:
 974

975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991

$$V(s) \leftarrow V(s) + \alpha_V * \left(\frac{R_{ABS} + Q(s, u)}{2} - V(s) \right)$$

When the feedback is complete, $Q(s, u)$ (the Q-value of the unchosen option) is replaced by the outcome of the unchosen option. α_V is an additional free parameter for the state value $V(s)$.

Concerning the learning phase, model simulation analysis (**Supp. Figure 3**) showed that, while the REFERENCE model matches the performance in the high-magnitude contexts in the learning phase ($\Delta EV=5$), it fails to capture the performance in low-magnitude contexts ($\Delta EV=0.5$). This is expected as the model does not implement range adaptation in any form. Concerning the transfer phase, the REFERENCE model reproduces a pattern that is qualitatively close to the observed results, but still less accurate compared to the RANGE model (out of sample likelihood comparison $LL_{RAN} = -96.79$ vs $LL_{REF} = -186.68$, $t(799) = 8.26$, $p < .0001$). To sum up, the REFERENCE model is strongly rejected by the learning phase results (where it essentially behaves like to the ABSOLUTE model) and weakly rejected by the transfer phase results, where it manages to capture the overall pattern, but in a less accurate manner.



992
993
994
995
996
997

Supp. Figure 3: generative performance of the RANGE model (black dots) compared to the REFERENCE model (white dots). Black lines represent the empirical averages. Colored squares indicate the s.e.m. around the empirical averages.

998

GLOBAL model

999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010

The RANGE model as we presented in the main text, does not contain any element to account for the block/interleaved effect. Here we propose a possible computational interpretation to account for the effects of this manipulation (more precisely the fact that contextual effects are exacerbated in block experiments). The key idea of this model is that the notion of ‘context’ can be break down into two components. The ‘local’ context is what we referred to as simply the context in the paper (essentially a pair of cues, or a state in the reinforcement learning framework). In addition to the local context we also postulate a ‘global’ context that integrate over a time scale larger than a trial (it could be understood as the current ‘value’ of the task). To instantiate these ideas, we built an alternative model (GLOBAL) that includes both “global” (or task-level) and local (or pair of options-level) contextual variables: $R_{MAX}(\text{task})$ and $R_{MAX}(\text{state})$. The $R_{MAX}(\text{state})$ is learnt similarly to the state-value in the REFERENCE model, except that it is not bounded to any particular pair of options:

1011

$$R_{\text{MAX}}(\text{task}) \leftarrow R_{\text{MAX}}(\text{task}) + \alpha_T * \left(\frac{R_{\text{ABS}} + Q(s, u)}{2} - R_{\text{MAX}}(\text{task}) \right)$$

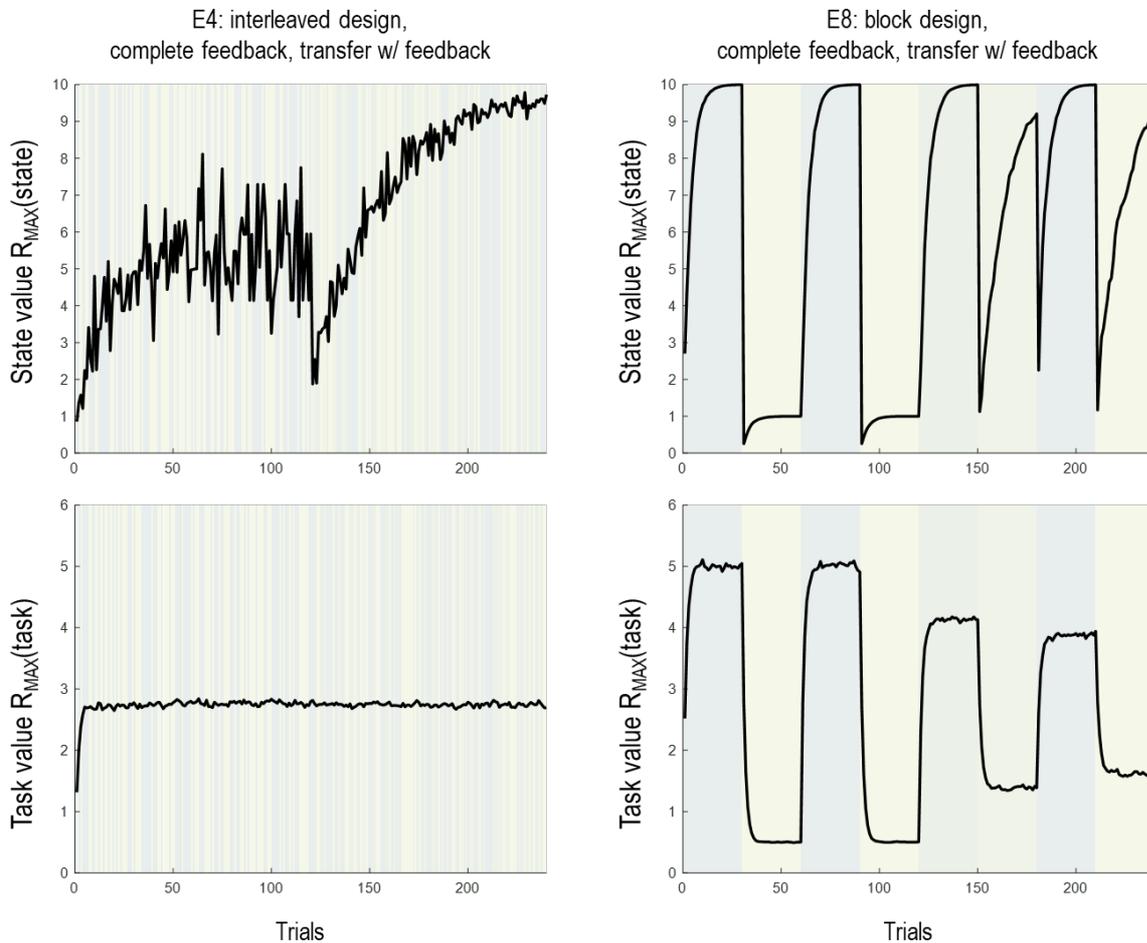
1012 When the feedback is complete, $Q(s, u)$ (the Q-value of the unchosen option) is replaced by the
 1013 outcome of the unchosen option. α_T is an additional free parameter for the $R_{\text{MAX}}(\text{task})$. The range
 1014 normalization rule (that we write here in its simplified manner that takes into account that $R_{\text{MIN}} =$
 1015 0 everywhere in our task) in the option value update rule of the GLOBAL model is as follows:
 1016

1017

$$Q(s, a) \leftarrow Q(s, a) + \alpha_Q * \left(\frac{R_{\text{ABS}}}{R_{\text{MAX}}(\text{state}) + R_{\text{MAX}}(\text{task}) + 1} - Q(s, a) \right)$$

1018

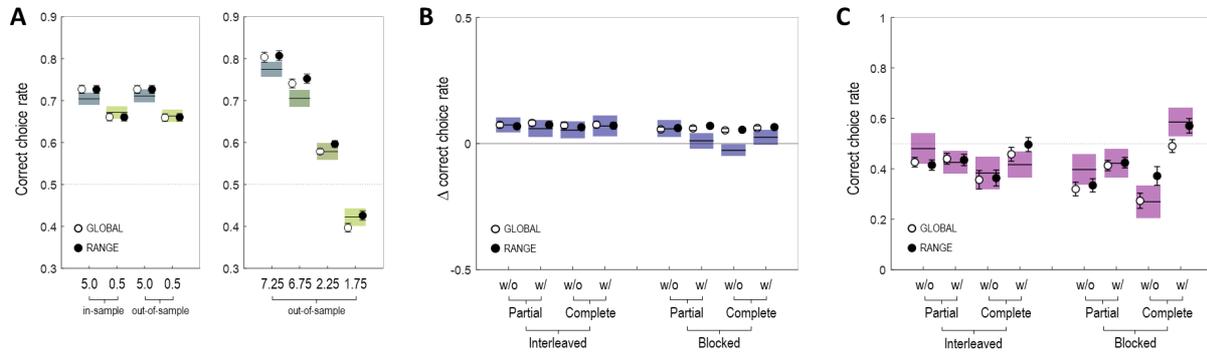
1019 This simple model accounts for increased contextual effects in block design, because in the block
 1020 design, $R_{\text{MAX}}(\text{task})$ and the $R_{\text{MAX}}(\text{state})$ remain coherent for longer time periods (**Supp. Figure**
 1021 **4**), thus allowing the summation of their effects. As shown in **Supp. Figure 5**, the model seems
 1022 qualitatively equal than the RANGE model, if not better at matching performance in most of the 8
 1023 different versions of the $\Delta\text{EV}=1.75$ context.
 1024
 1025



1026

1027 **Supp. Figure 4:** The figure illustrates the evolution across the experiment of the hidden variables
 1028 $R_{\text{MAX}}(\text{state})$ and $R_{\text{MAX}}(\text{task})$. Simulations concern E4 (interleaved design, complete feedback,

1029 transfer with feedback) and E8 (block design, complete feedback, transfer with feedback).
 1030 Background colors show the choice context (color code as in Figure 1).
 1031
 1032



1033
 1034
 1035 **Supp. Figure 5:** generative performance of the RANGE model (black dots) compared to the
 1036 GLOBAL model (white dots). Black lines represent the empirical averages. Colored squares
 1037 indicate the s.e.m. around the empirical averages.
 1038
 1039

1040 **REGRET model**

1041 Finally, we analyzed a model assuming that option values are purely encoded by outcome
 1042 comparison (akin to a relief/regret signal). A similar idea has been put forward by other studies
 1043 (7,35) where it proved successful in explain striatal neural activity and, to some extent, behavioral
 1044 data. This model has the strong handicap that it cannot be straight-forwardly extended to the partial
 1045 feedback case, where the outcome of the unchosen option is not showed. We therefore tested the
 1046 proposed model in the 4 experiments featuring complete feedback.

1047 Option values in the REGRET model are updated as follows, with R_C and R_U the outcomes of the
 1048 chosen option and unchosen option, respectively:

1049

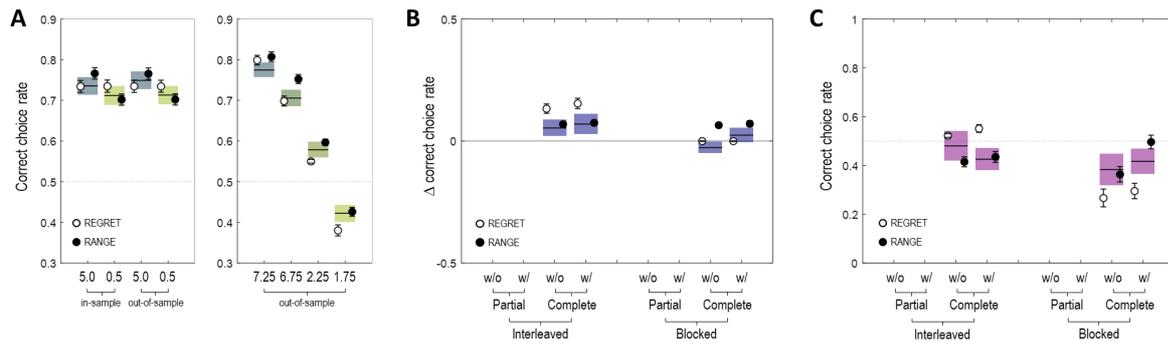
$$R_{\text{REG},t} = \begin{cases} 1 & \text{if } R_C > R_U \\ 0 & \text{if } R_C = R_U \\ -1 & \text{if } R_C < R_U \end{cases}$$

1050

$$Q_{t+1}(s, c) = Q_t(s, c) + \alpha_c * (R_{\text{REG},t} - Q_t(s, c))$$

$$Q_{t+1}(s, u) = Q_t(s, u) + \alpha_u * (R_{\text{REG},t} - Q_t(s, u))$$

1051
 1052
 1053
 1054 As clearly illustrated by the model simulations (**Supp. Figure 6**), the REGRET model does not fit
 1055 well the behavioral data, especially in the transfer phase, where it overestimates value inversion in
 1056 the $\Delta EV=1.75$ context. In other terms through a different mechanism, the REGRET model suffers
 1057 from the same problem the BINARY model: they predict to much option value context dependence.
 1058
 1059



1060
 1061 **Supp. Figure 7:** generative performance of the RANGE model (black dots) compared to the
 1062 REGRET model (white dots). Black lines represent the empirical averages. Colored squares
 1063 indicate the s.e.m. around the empirical averages.
 1064
 1065

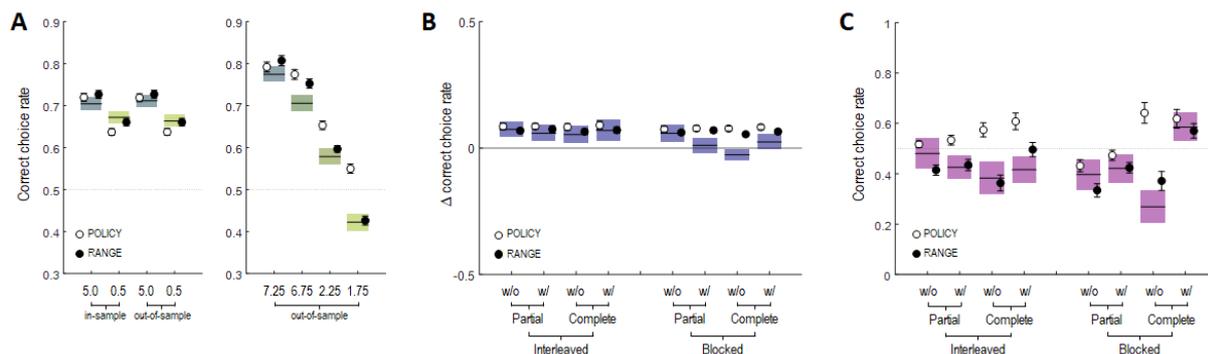
1066 POLICY model

1067 Finally, we considered a model that applies range normalization at the decision step, i.e., in the
 1068 softmax, instead of the outcome encoding stage as in the RANGE model. In this model (POLICY)
 1069 the probability of choosing option a over option b is defined by:
 1070

$$1071 P_t(s, a) = \frac{1}{1 + e^{\left(\beta^* \frac{Q_t(s,b) - Q_t(s,a)}{1 + \max\{Q_t(s,:)\} - \min\{Q_t(s,:)\}}\right)}}$$

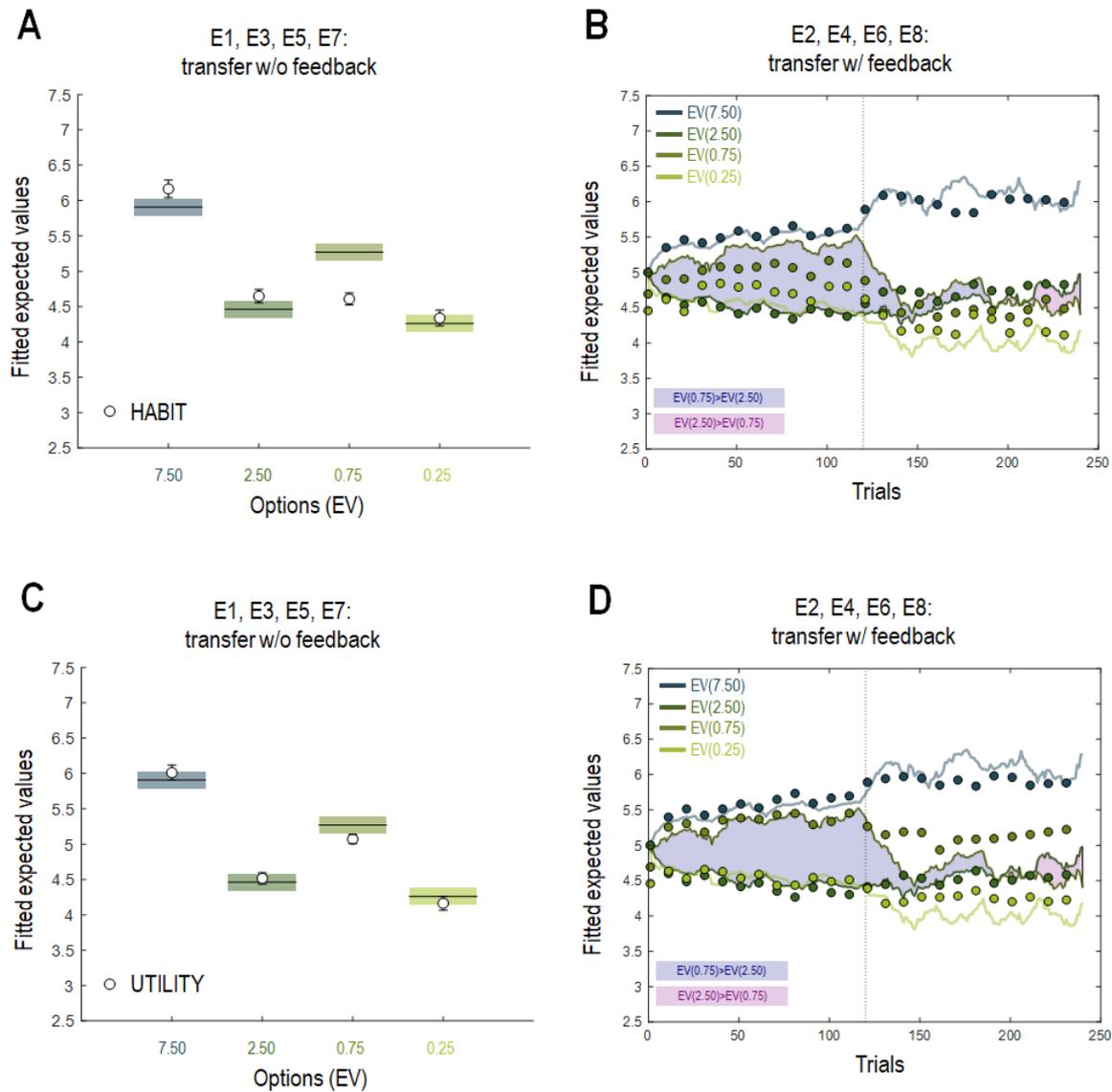
1072 Similarly to the RANGE model, the POLICY model is able to capture the magnitude difference in
 1073 the learning phase. In the transfer phase however, the POLICY model fails to predict the value
 1074 inversion in the $\Delta EV=1.75$ context. This is due to the fact that, despite the normalization process
 1075 within the softmax function, option values remain encoded in an absolute scale. Whereas in the
 1076 learning phase the POLICY model predicts a behavior compatible with the RANGE model, in the
 1077 transfer phase it predicts a behavior consistent with the ABSOLUTE model (**Supp. Fig. 6**).
 1078

1079



1080
 1081 **Supp. Figure 6:** generative performance of the RANGE model (black dots) compared to the
 1082 POLICY model (white dots). Black lines represent the empirical averages. Colored squares indicate
 1083 the s.e.m. around the empirical averages.
 1084
 1085

1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099



1100
1101
1102
1103
1104
1105
1106
1107

Supp. Figure 8. Ruling out habitual learning and marginally decreasing utility. (A-C) Average inferred option values for the behavioral data and simulated data for the experiments without trial-by-trial transfer feedback (white dots: HABIT (resp. UTILITY) model). (B-D) Trial-by-trial inferred option values for the behavioral data and simulated data for the experiments with trial-by-trial transfer feedback, where curves indicate trial-by-trial fit of each inferred option value, and colored dots indicate HABIT (resp. UTILITY) model simulations.

1108 **Supplementary references**

- 1109
- 1110 1. Louie K, Glimcher PW. Efficient coding and the neural representation of value. *Ann N Y Acad Sci.* 2012
1111 Mar;1251:13–32.
- 1112 2. Vlaev I, Chater N, Stewart N, Brown GDA. Does the brain calculate value? *Trends Cogn Sci.* 2011
1113 Nov;15(11):546–54.
- 1114 3. Cox KM, Kable JW. BOLD Subjective Value Signals Exhibit Robust Range Adaptation. *J Neurosci.* 2014 Dec
1115 3;34(49):16533–43.
- 1116 4. Nieuwenhuis S, Heslenfeld DJ, Alting von Geusau NJ, Mars RB, Holroyd CB, Yeung N. Activity in human
1117 reward-sensitive brain areas is strongly context dependent. *NeuroImage.* 2005 May 1;25(4):1302–9.
- 1118 5. Elliott R, Agnew Z, Deakin JFW. Medial orbitofrontal cortex codes relative rather than absolute value of
1119 financial rewards in humans. *Eur J Neurosci.* 2008 May;27(9):2213–8.
- 1120 6. Bavard S, Lebreton M, Khamassi M, Coricelli G, Palminteri S. Reference-point centering and range-
1121 adaptation enhance human reinforcement learning at the cost of irrational preferences. *Nat Commun.* 2018
1122 29;9(1):4503.
- 1123 7. Klein TA, Ullsperger M, Jocham G. Learning relative values in the striatum induces violations of normative
1124 decision making. *Nat Commun.* 2017 Jun 20;8(1):16033.
- 1125 8. Palminteri S, Khamassi M, Joffily M, Coricelli G. Contextual modulation of value signals in reward and
1126 punishment learning. *Nat Commun.* 2015 Aug 25;6:8096.
- 1127 9. Freidin E, Kacelnik A. Rational Choice, Context Dependence, and the Value of Information in European
1128 Starlings (*Sturnus vulgaris*). *Science.* 2011 Nov 18;334(6058):1000–2.
- 1129 10. Pompilio L, Kacelnik A. Context-dependent utility overrides absolute memory as a determinant of choice.
1130 *Proc Natl Acad Sci.* 2010 Jan 5;107(1):508–12.
- 1131 11. Rustichini A, Conen KE, Cai X, Padoa-Schioppa C. Optimal coding and neuronal adaptation in economic
1132 decisions. *Nat Commun.* 2017 Oct 31;8(1):1208.
- 1133 12. Webb R, Glimcher PW, Louie K. The Normalization of Consumer Valuations: Context-Dependent
1134 Preferences From Neurobiological Constraints. *Manag Sci* [Internet]. 2020 May 27 [cited 2020 Jul 27];
1135 Available from: <https://pubsonline.informs.org/doi/abs/10.1287/mnsc.2019.3536>
- 1136 13. Fontanesi L, Palminteri S, Lebreton M. Decomposing the effects of context valence and feedback information
1137 on speed and accuracy during reinforcement learning: a meta-analytical approach using diffusion decision
1138 modeling. *Cogn Affect Behav Neurosci.* 2019 Jun 1;19(3):490–502.
- 1139 14. Collins AGE, Frank MJ. How much of reinforcement learning is working memory, not reinforcement
1140 learning? A behavioral, computational, and neurogenetic analysis. *Eur J Neurosci.* 2012 Apr;35(7):1024–35.
- 1141 15. Rabin M. Diminishing Marginal Utility of Wealth Cannot Explain Risk Aversion. 2000 Jun 9 [cited 2020 Jul
1142 27]; Available from: <https://escholarship.org/uc/item/61d7b4pg>
- 1143 16. Miller KJ, Shenhav A, Ludvig EA. Habits without values. *Psychol Rev.* 2019;126(2):292–311.
- 1144 17. Landry P, Webb R. Pairwise Normalization: A Neuroeconomic Theory of Multi-Attribute Choice [Internet].
1145 Rochester, NY: Social Science Research Network; 2019 Jan [cited 2020 Nov 30]. Report No.: ID 2963863.
1146 Available from: <https://papers.ssrn.com/abstract=2963863>

- 1147 18. Palminteri S, Wyart V, Koehlin E. The Importance of Falsification in Computational Cognitive Modeling. *Trends Cogn Sci.* 2017 Jun;21(6):425–33.
1148
- 1149 19. Katahira K. The statistical structures of reinforcement learning with asymmetric value updates. *J Math*
1150 *Psychol.* 2018 Dec 1;87:31–45.
- 1151 20. Louie K, Glimcher PW, Webb R. Adaptive neural coding: from biological to behavioral decision-making. *Curr*
1152 *Opin Behav Sci.* 2015 Oct 1;5:91–9.
- 1153 21. Dumbalska T, Li V, Tsetsos K, Summerfield C. A map of decoy influence in human multialternative choice.
1154 *Proc Natl Acad Sci.* 2020 Oct 6;117(40):25169–78.
- 1155 22. Daviet R, Webb R. A Double Decoy Experiment to Distinguish Theories of Dominance Effects [Internet].
1156 Rochester, NY: Social Science Research Network; 2019 Mar [cited 2020 Nov 30]. Report No.: ID 3374514.
1157 Available from: <https://papers.ssrn.com/abstract=3374514>
- 1158 23. Gluth S, Kern N, Kortmann M, Vitali CL. Value-based attention but not divisive normalization influences
1159 decisions with multiple alternatives. *Nat Hum Behav.* 2020 Jun;4(6):634–45.
- 1160 24. Goodwin PB. Habit and Hysteresis in Mode Choice. *Urban Stud.* 1977;14(1):95–8.
- 1161 25. Dickinson A, Weiskrantz L. Actions and habits: the development of behavioural autonomy. *Philos Trans R*
1162 *Soc Lond B Biol Sci.* 1985 Feb 13;308(1135):67–78.
- 1163 26. Lally P, Jaarsveld CHM van, Potts HWW, Wardle J. How are habits formed: Modelling habit formation in the
1164 real world. *Eur J Soc Psychol.* 2010;40(6):998–1009.
- 1165 27. Thrailkill EA, Trask S, Vidal P, Alcalá JA, Bouton ME. Stimulus Control of Actions and Habits: A Role for
1166 Reinforcer Predictability and Attention in the Development of Habitual Behavior. *J Exp Psychol Anim Learn*
1167 *Cogn.* 2018 Oct;44(4):370–84.
- 1168 28. Akaishi R, Umeda K, Nagase A, Sakai K. Autonomous mechanism of internal choice estimate underlies
1169 decision inertia. *Neuron.* 2014 Jan 8;81(1):195–206.
- 1170 29. Neumann J von, Morgenstern O. *Theory of Games and Economic Behavior.* Princeton University Press; 1953.
1171 774 p.
- 1172 30. Bernoulli D. Exposition of a New Theory on the Measurement of Risk. *Econometrica.* 1954;22(1):23–36.
- 1173 31. Markowitz H. The Utility of Wealth. *J Polit Econ.* 1952 Apr 1;60(2):151–8.
- 1174 32. Kahneman D, Tversky A. Subjective probability: A judgment of representativeness. *Cognit Psychol.* 1972 Jul
1175 1;3(3):430–54.
- 1176 33. Loomes G, Sugden R. Regret Theory: An Alternative Theory of Rational Choice Under Uncertainty. *Econ J.*
1177 1982;92(368):805–24.
- 1178 34. Dayan P, Abbott LF. *Theoretical Neuroscience: Computational and Mathematical Modeling of Neural*
1179 *Systems.* Massachusetts Institute of Technology Press; 2001. 460 p.
- 1180 35. Li J, Daw ND. Signals in human striatum are appropriate for policy update rather than value prediction. *J*
1181 *Neurosci Off J Soc Neurosci.* 2011 Apr 6;31(14):5504–11.
- 1182 36. Palminteri S, Pessiglione M. Chapter 23 - Opponent Brain Systems for Reward and Punishment Learning:
1183 Causal Evidence From Drug and Lesion Studies in Humans. In: Dreher J-C, Tremblay L, editors. *Decision*

- 1184 Neuroscience [Internet]. San Diego: Academic Press; 2017 [cited 2020 Nov 30]. p. 291–303. Available from:
1185 <http://www.sciencedirect.com/science/article/pii/B9780128053089000233>
- 1186 37. Lebreton M, Bacily K, Palminteri S, Engelmann JB. Contextual influence on confidence judgments in human
1187 reinforcement learning. *PLOS Comput Biol*. 2019 avr;15(4):e1006973.
- 1188 38. Burke CJ, Baddeley M, Tobler PN, Schultz W. Partial Adaptation of Obtained and Observed Value Signals
1189 Preserves Information about Gains and Losses. *J Neurosci*. 2016 Sep 28;36(39):10016–25.
- 1190 39. Pischedda D, Palminteri S, Coricelli G. The Effect of Counterfactual Information on Outcome Value Coding
1191 in Medial Prefrontal and Cingulate Cortex: From an Absolute to a Relative Neural Code. *J Neurosci*. 2020 Apr
1192 15;40(16):3268–77.
- 1193 40. Conen KE, Padoa-Schioppa C. Partial Adaptation to the Value Range in the Macaque Orbitofrontal Cortex. *J*
1194 *Neurosci*. 2019 May 1;39(18):3498–513.
- 1195 41. Padoa-Schioppa C, Rustichini A. Rational Attention and Adaptive Coding: A Puzzle and a Solution. *Am Econ*
1196 *Rev*. 2014 May;104(5):507–13.
- 1197 42. Gigerenzer G. The Bias Bias in Behavioral Economics. *Rev Behav Econ*. 2018 Dec 30;5(3–4):303–36.
- 1198 43. Haselton MG, Nettle D, Andrews PW. The Evolution of Cognitive Bias. In: *The Handbook of Evolutionary*
1199 *Psychology* [Internet]. John Wiley & Sons, Ltd; 2015 [cited 2020 Jul 27]. p. 724–46. Available from:
1200 <https://onlinelibrary.wiley.com/doi/abs/10.1002/9780470939376.ch25>
- 1201 44. Frank MJ, Seeberger LC, O'reilly RC. By carrot or by stick: cognitive reinforcement learning in parkinsonism.
1202 *Science*. 2004 Dec 10;306(5703):1940–3.
- 1203 45. Girden ER. ANOVA: Repeated Measures. SAGE; 1992. 88 p.
- 1204 46. Sutton RS, Barto AG. Reinforcement Learning - An Introduction [Internet]. Mit Press; 1998 [cited 2017 Jun
1205 8]. Available from: <http://gen.lib.rus.ec/book/index.php?md5=C5BC41EB7F60C05D37473B4A9AC77A2A>
- 1206 47. Rescorla RA, Wagner AR. A theory of Pavlovian conditioning: Variations in the effectiveness of
1207 reinforcement and nonreinforcement. *Class Cond II Curr Res Theory*. 1972;2:64–99.
- 1208 48. Hergueux J, Jacquemet N. Social preferences in the online laboratory: a randomized experiment. *Exp Econ*.
1209 2015 Jun 1;18(2):251–83.
- 1210 49. Kahneman D. Maps of Bounded Rationality: Psychology for Behavioral Economics. *Am Econ Rev*. 2003
1211 Dec;93(5):1449–75.
- 1212 50. Shavit T, Sonsino D, Benzion U. A comparative study of lotteries-evaluation in class and on the Web. *J Econ*
1213 *Psychol*. 2001 Aug 1;22(4):483–91.
- 1214 51. Schoeffler M, Stöter F-R, Bayerlein H, Edler B, Herre J. An Experiment about Estimating the Number of
1215 Instruments in Polyphonic Music: A Comparison Between Internet and Laboratory Results. In: *ISMIR*. 2013.
- 1216 52. Reinecke K, Gajos KZ. LabintheWild: Conducting Large-Scale Online Experiments With Uncompensated
1217 Samples. In: *Proceedings of the 18th ACM Conference on Computer Supported Cooperative Work & Social*
1218 *Computing* [Internet]. Vancouver, BC, Canada: Association for Computing Machinery; 2015 [cited 2020 Jul
1219 27]. p. 1364–1378. (CSCW '15). Available from: <https://doi.org/10.1145/2675133.2675246>

1220
1221