

# Learning in visual regions as support for the bias in future value-driven choice

Sara Jahfari<sup>\*1,2</sup>, Jan Theeuwes<sup>3</sup>, Tomas Knapen<sup>1,3</sup>

<sup>1</sup> Spinoza Centre for Neuroimaging, Royal Netherlands Academy of Arts and Sciences (KNAW), The Netherlands

<sup>2</sup> Department of Psychology, University of Amsterdam, The Netherlands

<sup>3</sup> Department of Experimental and Applied Psychology, Vrije Universiteit van Amsterdam, The Netherlands

## Abstract

Learning biases decision-making towards higher expected outcomes. Cognitive theories describe this through the tracking of value and outcome evaluations within striatum and prefrontal cortex. Decisions however first require processing of sensory input, and to-date, far less is known about the learning perception interplay. This fMRI study (N=43), relates visual BOLD responses to value-beliefs during choice, and, signed prediction errors after outcomes. To understand these relationships, which co-occurred in striatum, we next evaluated relevance with the prediction of future value-based choices, using a separate transfer-phase with learning already established. We decoded choice outcomes with a 69% accuracy with a machine learning algorithm that was given trial-by-trial BOLD from visual regions alongside traditional prefrontal, and striatal regions. Importantly, this decoding of value-driven choice outcomes again showed an important role for visual activity. These results raise the intriguing possibility that value learning in visual cortex is supportive for the striatal bias towards valued options.

**Keywords:** Reinforcement learning, perceptual learning, decoding, Bayesian hierarchical modelling, random forest machine learning

---

\*Corresponding author: sara.jahfari@gmail.com

## 25 Introduction

26 In decision-making, our value beliefs bias future choices. This bias is shaped by the outcomes of similar  
27 decisions made in the past where the action, or stimulus chosen, becomes associated with a positive or  
28 negative outcome ('value beliefs'). The evaluation of value after an outcome, or the comparison of value in  
29 decisions, is traditionally associated with activity in the prefrontal cortex and striatum<sup>1-7</sup>.

30 To underset the bias in action selection midbrain dopamine neurons are thought to send a teaching signal  
31 towards the striatum and prefrontal cortex after an outcome<sup>8-10</sup>. In the striatum, future actions are facilitated  
32 by bursts in dopamine after positive outcomes or discouraged by dopamine dips after negative outcomes. The  
33 dorsal and ventral parts of the striatum are known to receive differential, but also overlapping, inputs from  
34 midbrain neurons<sup>7,11</sup>. Ventral and dorsal striatum have also been ascribed a differential role during learning  
35 by reinforcement learning theories. Here, the ventral parts of the striatum are involved with the prediction  
36 of future outcomes through the processing of prediction errors, whereas the dorsal striatum uses the same  
37 information to maintain action values as a way to bias future actions towards the most favored option<sup>4,12,13</sup>.  
38 Intriguingly, however, before many of these value-based computations can take place, stimuli first have to  
39 be parsed from the natural world, an environment where most reward predicting events are perceptually  
40 complex. This suggests that sensory processing might be an important integral part of optimized value-based  
41 decision-making.

42 Here, we investigate whether choice outcomes can modulate the early sensory processing of perceptually  
43 complex stimuli to help bias future decisions. Recent neurophysiological studies find visually responsive  
44 neurons in the tail of the caudate nucleus, which is part of the dorsal striatum<sup>14,15</sup>. These neurons encode and  
45 differentiate stable reward values of visual objects to facilitate eye movements towards the most valued target,  
46 while at the same time inhibiting a movement towards the lesser valued object<sup>16</sup>. Critically, differential  
47 modulations are also observed in the primary visual cortex where stronger cortical responses are seen for  
48 objects with higher values<sup>17,18</sup>, which is consistent with the response of visual neurons in the caudate. As  
49 visual cortex is densely connected to the striatum<sup>19,20</sup>, prioritized visual processing of high-value stimuli  
50 could aid the integration of information regarding the most-valued choice in the striatum<sup>21-24</sup>. To understand  
51 these visual-striatal interactions, we focus on a more detailed parsing of the underlying computations.

52 Specifically, we explored two questions by reanalyzing fMRI data from a probabilistic reinforcement learning  
53 task using faces as visual stimuli<sup>25</sup> (Figure 1a). First, we focus on the interplay between learning and visual  
54 activity in the fusiform face area (FFA) and occipital cortex (OC). Here, with the use of a Bayesian hierarchical  
55 reinforcement learning model (Figure 1b) we outline how trial-by-trial estimates of action values ( $Q$ -value)  
56 and reward prediction errors (RPE) relate to the BOLD response of visual regions and the striatum<sup>26,27</sup>  
57 (Figure 1c). Second, we analyze data from a follow-up transfer phase, where the learning of value was already  
58 established. In our analysis, the importance of visual brain activity in the prediction, or decoding, of future  
59 value-based decisions is evaluated by using a supervised Random Forest (RF) machine learning algorithm<sup>28,29</sup>.  
60 Specifically, transfer phase single-trial BOLD estimates from anatomically defined visual, prefrontal, and  
61 subcortical regions are combined by RF to predict, or decode, choice outcomes in a separate validation  
62 set. We focus on classification accuracy, and the relative importance of each brain region in the correct  
63 classification of future value-based decisions.

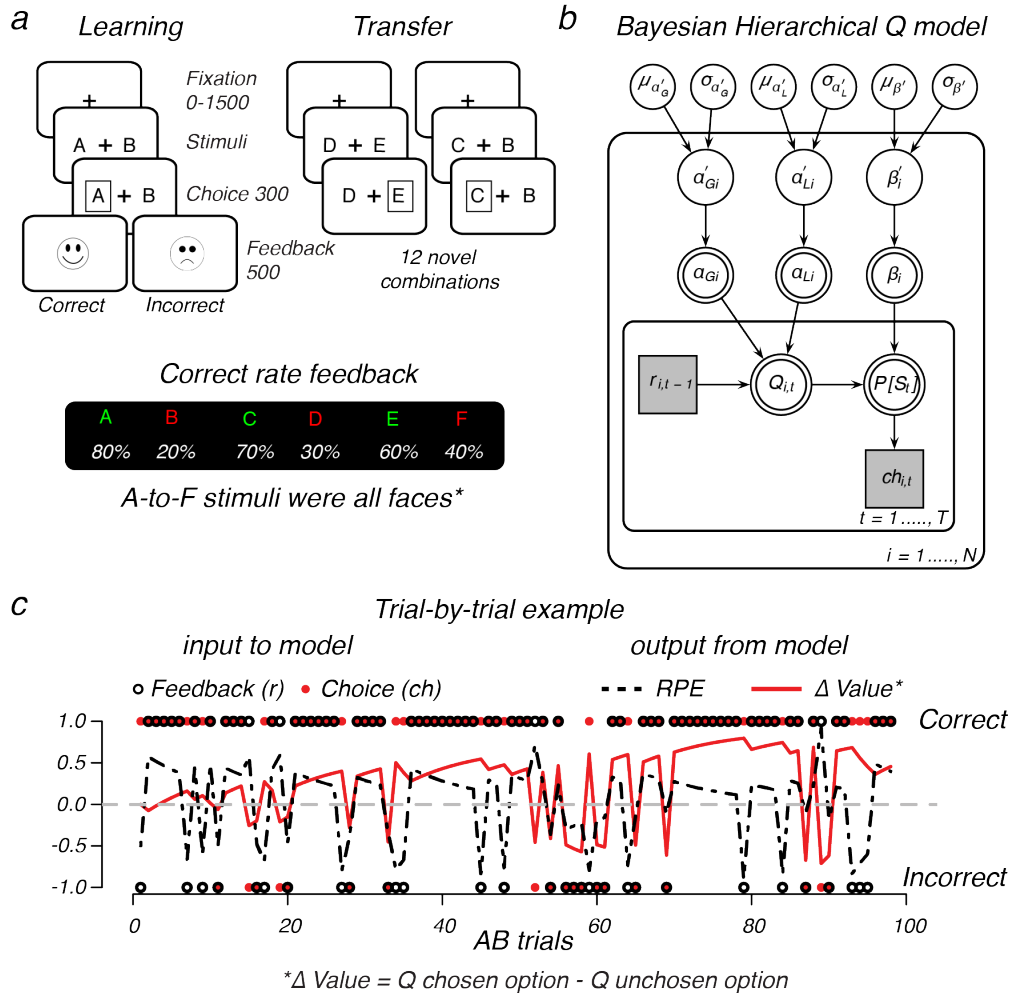


Figure 1: **Design and Model.** **a)** Reinforcement learning task using faces. During learning, two faces were presented on each trial, and participants learned to select the optimal face identity (A, C, E) through probabilistic feedback (% of correct is shown beneath each stimulus). The learning-phase contained three face pairs (AB, CD, ED) for which feedback was given. In a follow-up transfer phase these faces were rearranged into 12 novel combinations to assess learning. These trials were identical to learning trials, with the exception of feedback. \*Example faces were removed for the publication on BioRxiv, for an impression see 25, or the Radboud face database from where the faces were originally selected (<http://www.socsci.ru.nl:8180/RaFD2/RaFD>). **b)** Graphical  $Q$ -learning model with hierarchical Bayesian parameter estimation. The model consists of an outer subject ( $i = 1, \dots, N$ ), and an inner trial plane ( $t = 1, \dots, T$ ). Nodes represent variables of interest. Arrows are used to indicate dependencies between variables. Double borders indicate deterministic variables. Continuous variables are denoted with circular nodes, and discrete with square nodes. Observed variables are shaded in grey (see methods for details about the fitting procedure). **c)** Illustration of the observed trial-by-trial input (i.e., the choice made, and feedback received), and output (i.e.,  $Q$  for the chosen and unchosen stimulus,  $\Delta$ Value, and RPE) of the model given the estimated variability in learning rates from either positive ( $\alpha_{Gi}$ ) or negative ( $\alpha_{Li}$ ) feedback, and the tendency to exploit  $\beta$  higher values  $i$ .

## 64 Results

65 To understand how value learning relates to the activity pattern in perceptual regions we reanalyzed the  
66 behavioral and fMRI recordings of a recent study<sup>25</sup>. In this study, BOLD signals were recorded while  
67 participants performed a reinforcement learning task using male or female faces, and a stop-signal task (which  
68 was discussed in 25). The fusiform face area (FFA) was localized using a separate experimental run. 49  
69 young adults (25 male; mean age = 22 years; range 19-29 years, 43 analyzed, see Methods) participated in  
70 the study. As shown in Figure 1a, in the reinforcement learning task participants learned to select among  
71 choices with different probabilities of reinforcement (i.e., AB 80:20, CD 70:30, and EF 60:40). A subsequent  
72 transfer phase, where feedback was omitted, required participants to select the optimal option among novel  
73 pair combinations of the faces that were used during the learning phase (Figure 1a).

## 74 Model and Behavior

75 In the learning phase, subjects reliably learned to choose the most optimal face option in all pairs. For each  
76 pair the probability of choosing the better option was above chance ( $p$ 's < .001), and the effect of learning  
77 decreased from AB (80:20) and CD (70:30) to the most uncertain EF (60:40) pair ( $F(2, 84) = 13.74, p < .0001$ ).  
78 At the end of learning, value beliefs differentiating the optimal (A, C, E) from the sub-optimal (B, D, F) action  
79 were very distinct for the AB and CD face pairs but decreased with uncertainty ( $F(2, 84) = 39.70, p < 0.0001$ ,  
80 Figure 2a). Value beliefs were estimated using the individual subject parameters of the  $Q$ -learning model  
81 that best captured the observed data (Figure 2b-e; reproduced from 25 to show performance).

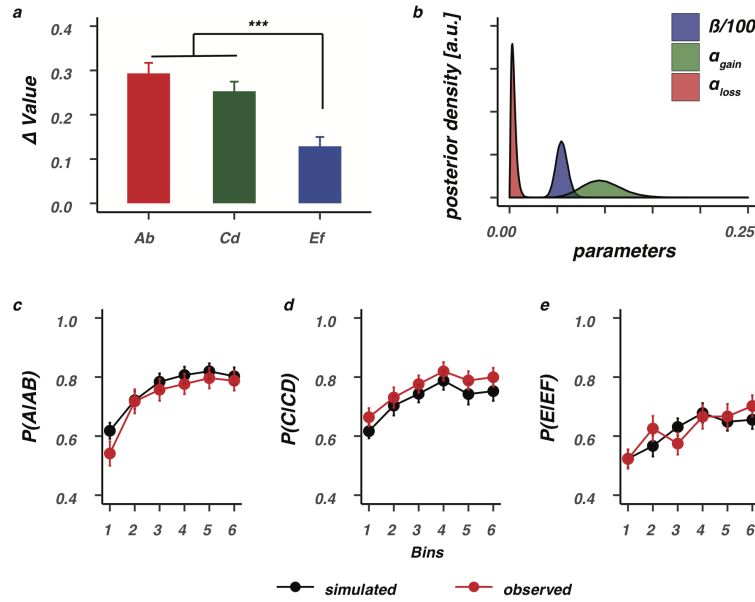


Figure 2: **Value differentiation and model performance.** a) Value differentiation ( $\Delta$ Value) for the selection of the optimal (A,C,E) stimuli over the suboptimal (B,D,F) stimuli decreased as a function of feedback reliability, and was smallest for the most uncertain EF stimuli. \* \* \* =  $p < 0.0001$ , Bonferroni corrected. b) Group-level posteriors for all Q-learning parameters. The bottom row shows model performance, where data was simulated with the estimated individual subject parameters and evaluated against the observed data for the AB (c), CD (d), or EF (e) pairs. Bins contain  $\pm 16$  trials. Error bars represent standard error of the mean (SEM).

## 82 BOLD is modulated by reliable value differences between faces in striatal and 83 visual regions

84 For each pair of faces presented during the learning phase (AB, CD, EF) we asked how the BOLD signal  
85 time-course in striatal and visual regions relates to trial-by-trial value beliefs about the two faces presented  
86 as a choice. First, as a reference, we focused on the activity pattern of three striatal regions. Results showed  
87 BOLD responses in dorsal (caudate, putamen) but not ventral (accumbens) striatum to be differentially  
88 modulated by the estimated value beliefs of the chosen face ( $Q_{chosen}$ ), in comparison to value beliefs about  
89 the face that was not chosen ( $Q_{unchosen}$ ). Thus, BOLD responses in the dorsal striatum were modulated  
90 more strongly by value beliefs about the chosen stimulus ( $Q_{chosen}$ ; Figure 3a bottom row). Critically, this  
91 differential modulation was only observed with the presentation of AB faces where value differences were  
92 most distinct because of the reliable feedback scheme. Next, we evaluated the relationship between value and  
93 BOLD in the FFA, and OC. Again, only with the presentation of the AB face option, trial-by-trial BOLD  
94 fluctuations were differentially modulated by values of the chosen versus not chosen face option (Figure  
95 3b bottom row). These evaluations highlight how the BOLD response in striatal and perceptual regions is  
96 especially sensitive to values of the (to-be) chosen stimulus when belief representations are stable and distinct.

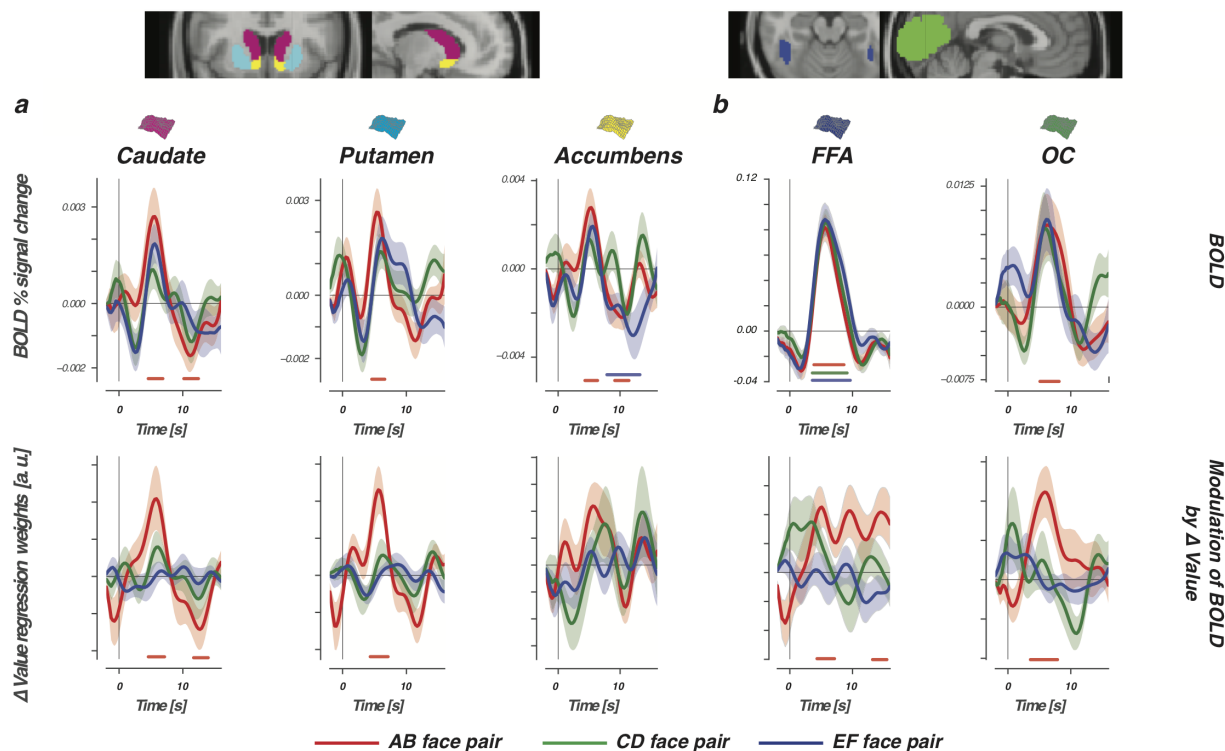


Figure 3: **BOLD and the modulation of  $\Delta$ Value in the learning phase.** Top row shows the BOLD signal time course, time-locked to presentations of AB (80:20, red lines), CD (70:30, green lines), and EF (60:40, blue lines) face pairs, for three striatal regions (a) and two perceptual regions (b). Bottom row displays differential modulation by value ( $\Delta$ Value = modulation  $Q_{chosen}$  – modulation  $Q_{unchosen}$ ). Horizontal lines show the interval in which modulation was significantly stronger for  $Q_{chosen}$ . With the presentation of AB faces, BOLD responses in the dorsal striatum (caudate and putamen) and visual regions (FFA and OC) were modulated more by values of the chosen stimulus when compared to values of the unchosen stimulus. Differential AB value modulation was not significant in the ventral striatum (i.e., accumbens). Nor did we observe any differential value modulations with the presentation of the more uncertain CD and EF pairs. Confidence intervals were estimated using bootstrap analysis across participants ( $n = 1000$ ), where the shaded region represents the standard error of the mean across participants (bootstrapped 68% confidence interval).

## 97 Reward prediction errors in striatal and visual regions

98 Our findings so far described relationships between BOLD and value time-locked to the moment of stimulus  
 99 presentation – i.e., when a choice is requested. Learning occurs when an outcome is different from what  
 100 was expected. We therefore next focused on modulations of the BOLD response when participants received  
 101 feedback. Learning modulations were explored by asking how trial-by-trial BOLD responses in perceptual and  
 102 striatal regions relate to either signed (outcome was better or worse than expected) or unsigned (magnitude of  
 103 expected violation) reward prediction errors<sup>30</sup>. Consistent with the literature, BOLD responses in all striatal  
 104 regions were modulated by signed RPEs, with larger responses after positive RPEs or smaller responses after  
 105 negative RPEs (Figure 4a bottom row). Activity in the accumbens (ventral striatum) was additionally tied  
 106 to unsigned RPEs in the tail of the BOLD time-course, with larger violations (either positive or negative)  
 107 tied to smaller dips. Consistently, estimated BOLD responses in both visual regions were modulated by the

108 signed RPE, and once more mirrored the striatal modulations with stronger positive RPEs eliciting stronger  
 109 BOLD responses (Figure 4b bottom row). FFA BOLD responses were additionally modulated by unsigned  
 110 RPEs. However, in contrast to the relationship found between unsigned RPEs and the accumbens, the FFA  
 111 modulation was positive and co-occurred with the modulation of the signed RPE. That is, bigger violations  
 112 and more positive outcomes each elicited a stronger response in the FFA.

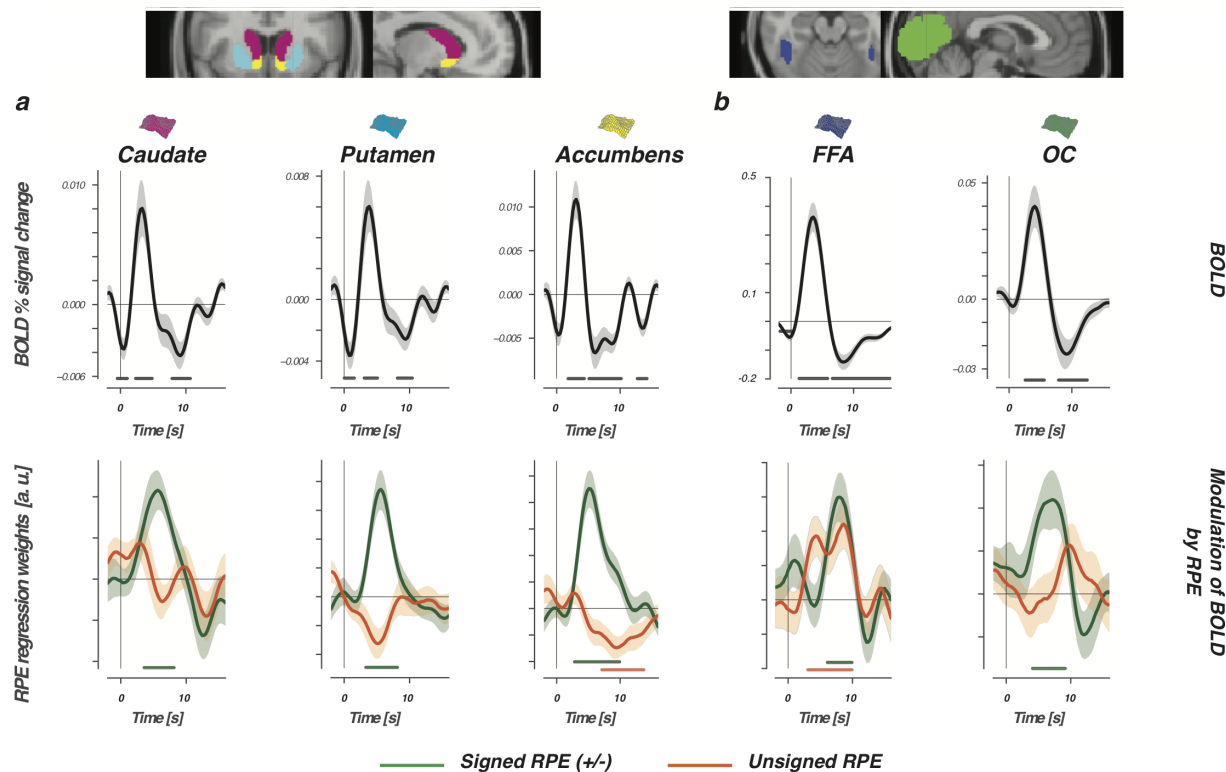


Figure 4: **Reward prediction errors modulate BOLD in striatal and visual regions.** The top row shows the FIR-estimated BOLD signal time-course, which was time-locked to the presentation of choice feedback and evaluated for three striatal regions (a) and two perceptual regions (b). Bottom row displays modulations of the estimated BOLD time-course by signed (green lines), or unsigned (orange lines) RPEs. The horizontal lines represent the interval in which signed or unsigned RPEs contributed significantly to the modulation of BOLD in the multiple regression. Note that both variables were always evaluated simultaneously in one GLM.

### 113 Can past learning in visual regions support the prediction of future value-based 114 decisions?

115 Stable value representations and reward prediction errors both modulated the activity of visual and striatal  
 116 regions. These modulations in the striatum are described to bias future actions towards the most favored  
 117 option (the dorsal striatum), or to predict future reward outcomes (the ventral striatum). To better understand  
 118 the value and RPE modulations observed in visual regions, we next assessed the importance of these visual  
 119 regions alongside the striatum in the correct classification (decoding) of future value-driven choice outcomes.  
 120 Here, activity of prefrontal regions was added to the importance evaluation based on our previous work with



121 this data in the transfer phase<sup>25</sup> (please see supplementary Figures 1&2 for the evaluation of these regions  
 122 during learning).

123 In the transfer phase, participants had to make a value-driven choice based on what was learned before, i.e.,  
 124 during the learning phase. To specify the relevance of visual regions in the resolve of value-driven choice  
 125 outcomes, in the transfer phase, a random forest (RF) classifier was used<sup>28,29</sup> (Please see Figure 5a-c for the  
 126 procedure). The RF classifier relies on an ensemble of decision trees as base learners, where the prediction of  
 127 each trial outcome is obtained by a majority vote that combines the prediction of all decision trees (Figure 6a).  
 128 To achieve controlled variation, each decision tree is trained on a random subset of the variables (i.e. subset of  
 129 columns shown in Figure 5a), and a bootstrapped sample of data points (i.e. trials). Importantly, we ensured  
 130 that the forest was not simply learning the proportion of optimal choices in the transfer phase by training all  
 131 models on balanced draws from the training set with equal numbers of optimal and sub-optimal choices.

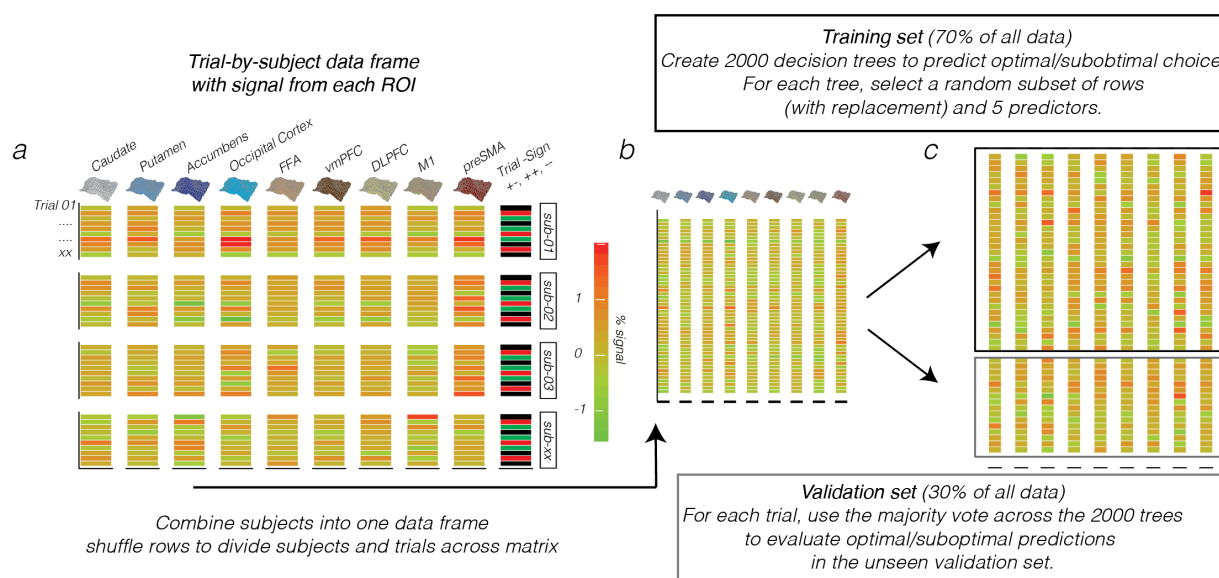


Figure 5: **Random Forest input and data-structure.** **a)** Trial-by-subject data matrix with the % signal change drawn for each choice trial in the transfer-phase (rows) from 9 a-priori defined regions of interest (columns). In addition to the ROI data, the matrix contained a column with the identity of participants (sub-01, etc) and Trial Sign, which specified a choice between two positives (+/+; AC, AE, CE), negatives (-/-, BD, BF, DF), or between a negative and positive option (+/-, e.g., AD, CF, etc) given the feedback scheme in the learning-phase. **b)** The individual subject data frames were then combined into one matrix, in which the rows were subsequently shuffled to randomly distribute trials and subjects across the rows. **c)** This matrix was then divided into a training set (2/3 of the data) for the creation of 2000 decision trees of which the majority vote on each trial is then used to evaluate the predictive accuracy of optimal/suboptimal choices in a separate validation set (1/3 of the data).

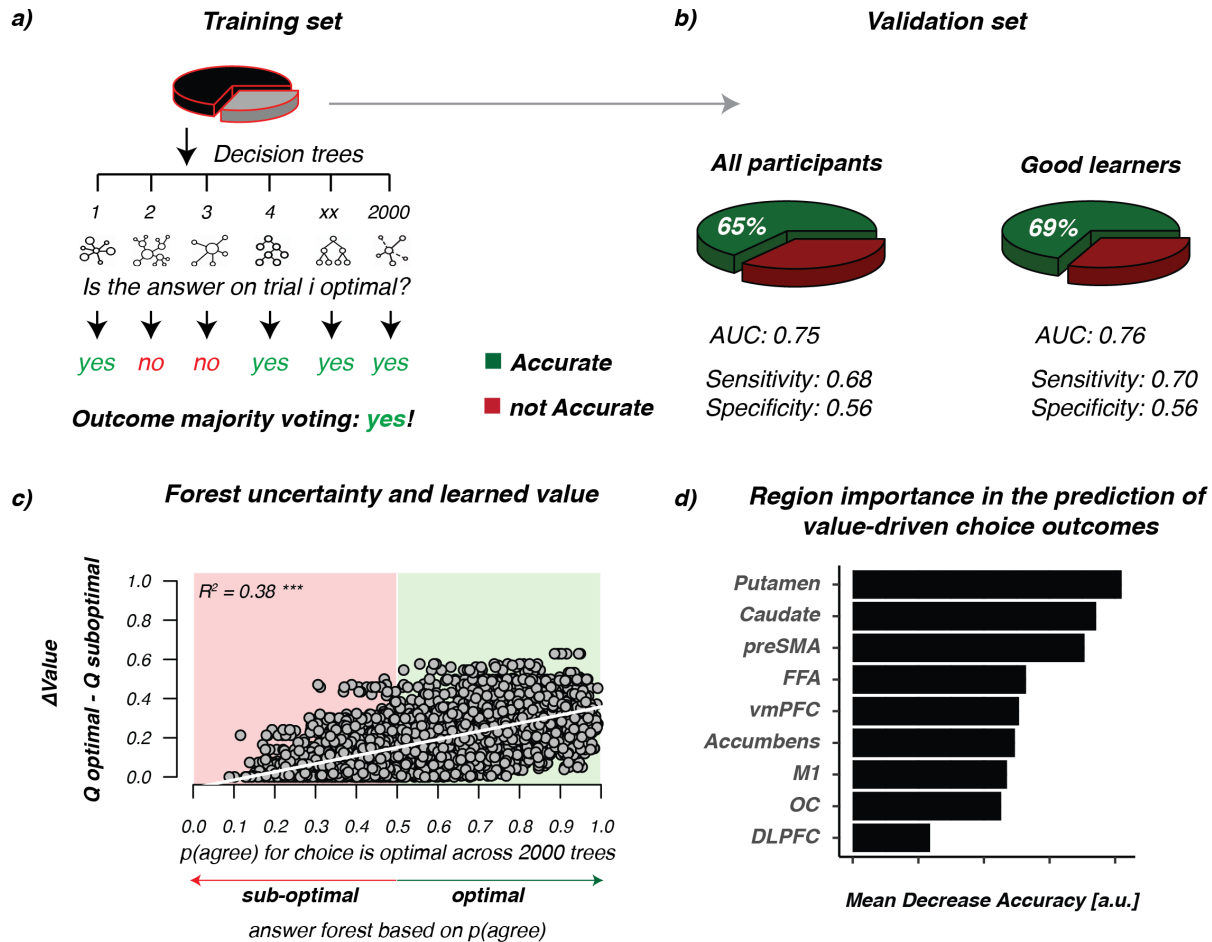
132 Evaluation of all participants resulted in a classification accuracy of 65% ( $AUC = 0.75$ ) using the trial-by-trial  
 133 BOLD estimates from the ROIs and increased to 69% with the evaluation of the good learners ( $AUC = 0.76$ ;  
 134  $N = 34$ , criteria: accuracy > 60% across all three learning pairs). Hence, in 65 (all participants) or 69 (good  
 135 learners) out of 100 trials the forest correctly classified whether participants would pick the option with  
 136 the highest value (optimal choice) or not (sub-optimal choice) in the validation set. RF predictions were  
 137 substantially lower when labels of the validation set were randomly shuffled (accuracy: all participants= 52%;



138 good learners= 56%).

139 The improvement of accuracy with the evaluation of only the good learners is remarkable because the classifier  
140 was given less data to learn the correct labelling (fewer subjects/trials) and implied that the 2000 decision  
141 trees were picking up information related to past learning. Further support for this important observation  
142 was found by asking how the uncertainty of each prediction (defined as the proportion of agreement in the  
143 predicted outcome among the 2000 trees for each trial) relates to the difference in value beliefs ( $\Delta$ Value)  
144 about the two options presented on each trial (computed using the end  $Q_{beliefs}$  of participants at the end  
145 of learning about face A-to-E). As plotted in Figure 6c, the uncertainty in predicting that a trial choice  
146 outcome is optimal – defined as the proportion of disagreement among the 2000 decision trees - decreased  
147 with larger belief differences in the assigned values (please see supplementary Figure 3 for the evaluation of  
148 all participants).

149 Besides providing insights into how BOLD responses in the transfer-phase contribute to predict value-driven  
150 choice outcomes (i.e., whether participants would choose the option with the highest value given past learning)  
151 the RF algorithm additionally outputs a hierarchy, thereby ranking the contribution of each region in the  
152 achieved classification accuracy. Figure 6d shows the ranking of all ROIs for good learners where the model  
153 had the highest predictive accuracy. First, regions in the dorsal striatum were most important, which aligned  
154 well with both the literature and the BOLD modulations we found by  $\Delta$ Value and RPE during the learning  
155 phase. These regions were next followed by the preSMA. Evaluation of this region during the learning phase  
156 showed no modulations by  $\Delta$ Value or RPE on BOLD ( supplementary Figure 1&2). Nevertheless, this region  
157 is typically associated with choice difficulty/conflict and might be essential in the resolve of a choice when  
158 value differences are small. Remarkably, the third region in this hierarchy was the FFA. In a task where  
159 participants pick the most valued face based on past learning, this ranking of the FFA just above the vmPFC  
160 and accumbens (ventral striatum) implies that the  $\Delta$ Value and RPE modulations of BOLD observed during  
161 learning could function to strengthen the recognition of valuable features. Note, however, that with the  
162 evaluation of all participants – including some who were less good in learning – the ranking of both the FFA  
163 and vmPFC was much lower (please see supplementary Figure 3b), which might be caused by more noise  
164 across the group in learning. We will return to this point in the discussion.



**Figure 6: Random Forest performance and importance ranking.** The prediction of value-driven choice outcomes in the transfer phase using trial-by-trial BOLD responses from striatal, perceptual, and prefrontal cortex regions. **a)** Overview of the Random Forest approach where the training-set is used to predict choice outcomes for each trial by using the majority vote of 2000 different decision trees. Each tree is built using a different set, or sample, of trials and predictors from the training set. The forest is trained on a training set sampled from all participants ( $N=43$ ), or only ‘the good learners’ ( $N=34$ ). **b)** Shows the classification, or decoding, accuracy (green) given the separate unseen validation sets, for all participants and good learners. **c)** Plotted relationship between forest uncertainty (i.e., proportion of agreement across 2000 trees), on each prediction/trial (x-axis) and  $\Delta$ Value (y-axis) for the model with the highest accuracy (i.e., the good learners).  $\Delta$ Value was computed for each trial in the transfer phase by using the end beliefs ( $Q$ ) that participants had about each stimulus (A-to-F) at the end of the learning phase. Forest uncertainty is defined as the proportion of trees saying ‘yes! the choice on this trial was optimal’. When this ratio is below 0.5 the forest will predict ‘no’ (sub-optimal), otherwise the prediction is ‘yes! the choice on this trial was optimal’ (optimal).  $R^2$ =adjusted  $R^2$ . Note that, the same pattern was found for all participants ( $R^2=0.41^{***}$ , please see supplementary Figure 3). **d)** Plotted ranking of the ROIs in their contribution to the predictive accuracy of the best performing model (i.e., good learners).

## 165 Discussion

166 This study provides novel insights into how reinforcements modulate visual activity and specifies its potential  
167 in the prediction of future value-driven choice outcomes. First, by focusing on how participants learn, we find  
168 BOLD in visual regions to change with trial-by-trial adaptations in value beliefs about the faces presented,  
169 and then to be subsequently scaled by the signed RPE after feedback. Next, the relevance of these observed  
170 value and feedback modulations was sought by exploring the prediction of future value-driven choice outcomes  
171 in a follow-up transfer phase where feedback was omitted. Our machine learning algorithm here shows a  
172 classification accuracy of 69% for participants who were efficient in learning by combining trial-by-trial BOLD  
173 estimates from perceptual, striatal, and prefrontal regions. The evaluation of region importance in these  
174 predictions ranked the FFA just after the dorsal striatum and the preSMA, thereby showing an important  
175 role for visual regions in the prediction of future value-driven choice outcomes in a phase where learning is  
176 established.

177 In a choice between two faces, BOLD responses in both the dorsal striatum and perceptual regions were  
178 affected more by values of the chosen face, relative to the unchosen face. Across three levels of uncertainty,  
179 we only observed the differential modulation of value on BOLD when belief representations were stable.  
180 This specificity aligns with neuronal responses to perceptual stimuli in the caudate tail<sup>16</sup>, visual cortex<sup>31,32</sup>,  
181 and imaging work across sensory modalities<sup>17,18,33-35</sup>, where it fuels theories in which the learning of stable  
182 reward expectations can develop to modulate, or sharpen, the representation of sensory information critical  
183 for perceptual decision making<sup>31,35</sup>.

184 After a choice was made, feedback modulations of signed ('valence') and unsigned ('surprise') RPEs<sup>30</sup> were  
185 evaluated on BOLD responses, by using an orthogonal design where the unsigned and signed RPE compete to  
186 explain BOLD variances. Both visual and striatal regions respond to prediction errors<sup>36</sup>. In the striatum both  
187 valence and surprise are thought to optimize future action selection in the dorsal striatum, or the prediction of  
188 future rewards in the ventral striatum. In perceptual regions, a mismatch between the expected and received  
189 outcome is often explained as surprise where a boost in attention or salience changes the representation of  
190 an image without a representation of value per se. We found positive modulatory effects of signed RPEs  
191 in all striatal regions, as well as, in the FFA and OC. Concurrently, modulations of unsigned RPEs were  
192 only observed in the accumbens (ventral striatum) and FFA, where notably the direction of modulation  
193 was reversed. We speculate that this contrast arises from the differential role of the regions. In the FFA,  
194 specialized and dedicated information processing is essential to quickly recognize valuable face features.  
195 Complementary boosts of surprise and valence here could prioritize attention towards the most rewarding  
196 face feature to strengthen the reward association in memory, or help speed up future recognition<sup>37-39</sup>. In  
197 the accumbens, boosted effects of positive valence on BOLD were dampened by larger mismatches. Large  
198 mismatches in what was expected are rare in stable environments. We therefore reason that in the accumbens  
199 the contrast between valence and surprise could function as a scale to refine learning, eventually leading to  
200 more reliable predictions of future rewards.

201 Whereas BOLD in the ventral striatum was shaped by both signed and unsigned RPEs, the dorsal striatum  
202 was sensitive to differential value up-to a choice and signed RPEs with the presentation of feedback<sup>34,40-43</sup>.  
203 The concurrent modulation of differential value in the primary motor cortex (please see M1 in supplementary  
204 Figure 1) associates the dorsal striatum with the integration of sensory information<sup>16,44-46</sup>, where increased  
205 visual cortex BOLD responses to faces with the highest value could potentially help bias the outcome of a

206 value-driven choice.

207 We explored this line of reasoning with the prediction of value-driven choice outcomes in a follow-up transfer  
208 phase after leaning. In recent years, machine learning approaches have become increasingly important in  
209 neuroscience<sup>47-50</sup>, where the ease of interpretation has often motivated a choice for linear methods above  
210 non-linear methods<sup>49,51</sup>. Despite the latter being less constrained and able to reach a better classification  
211 accuracy by capturing non-arbitrary, or unexpected relationships<sup>52</sup>. Value-driven choices after a phase of  
212 initial learning are influenced by the consistency of past learning, memory updating, and attention. All  
213 of these processes are affected by both linear and non-linear neurotransmitter modulations<sup>53-56</sup>. Our RF  
214 approach was unconstrained by linearity with classification accuracies well above chance and improved with  
215 the evaluation of only the good learners; despite substantial decreases in data given to the algorithm to  
216 learn the correct labelling. Critically, we additionally found that the uncertainty of trial-by-trial predictions  
217 made by RF is tied to the differentiability of value beliefs – an index that we could compute for the novel  
218 pair combination in the transfer phase by using the value ( $Q$ ) beliefs that participants had about each face  
219 at the end of learning. These results showcase how trial-by-trial BOLD fluctuations in striatal, prefrontal,  
220 and sensory regions can be combined by machine learning, or decoding, algorithms to reliably predict the  
221 outcome of a value-driven choice. Where we refine the interpretation of non-linear predictions by combining  
222 the RF output with cognitive computational modelling. With this combination we essentially show how the  
223 uncertainty of RF predictions is tied to value beliefs acquired with learning in the past.

224 An important evaluation intended with our machine learning approach was the ranking of regions by their  
225 contribution to the predictive (decoding) accuracy in the transfer phase. After the observed modulations of  
226 BOLD in the learning phase this explorative analysis sought the relevance of learning-BOLD relationships in  
227 the resolve of future choices. Here, the ranking made by RF first identified signals from the dorsal striatum  
228 (putamen and caudate) as most important followed by the preSMA, and then most notably, visual regions.  
229 That is, when the quality of leaning was high across participants, FFA ranked just above traditional regions  
230 such as the vmPFC and the accumbens<sup>2,5,6,57</sup>. Notably, FFA was replaced by OC in ranking with the  
231 evaluation of all participants (please see supplementary Figure 3b). This difference could occur because  
232 the quality of learning was more variable across all participants, or because RF predictions based on the  
233 heterogeneous data from all participants were less accurate. In general, the shift in ranking implies that when  
234 learning is less consistent choice outcomes are better predicted by fluctuations in OC - perhaps with the  
235 identification of rewarding low-level features. With better or more consistent learning, however, participants  
236 should increasingly rely on memory and specialized visual areas. Thus, search for specific face features  
237 associated with high value by recruiting the FFA in the visual ventral stream. Consistent with this reasoning  
238 recent neuronal recordings show rapid visual processing of category-specific value cues in the ventral visual  
239 stream. These specific value cues are only seen for well-learned reward categories, and critically, precede the  
240 processing of value in prefrontal cortex<sup>59</sup>.

241 We note that although BOLD fluctuations in the preSMA ranked second in the prediction of value-driven  
242 choice outcomes, no reliable modulations of BOLD were observed by either differential value or RPEs in the  
243 learning phase. The preSMA is densely connected to the dorsal striatum and consistently associated with  
244 action-reward learning<sup>60</sup>, or choice difficulty<sup>61</sup>. The lack of associations in this study might result from our  
245 noisier estimates of the BOLD response that is typical for regions in the prefrontal cortex<sup>62,63</sup>, the anatomical  
246 masks selected, or smaller variability across trials in the learning phase (i.e., 3 pairs in learning-phase vs 15  
247 pairs in transfer-phase). Nevertheless, the importance indicated by RF, combined with our previous analysis

248 of this transfer phase data<sup>25</sup>, implies an important role for the preSMA in the resolve of value-driven choices  
249 in concert with the striatum. More research with optimized sequences to estimate BOLD in PFC is required  
250 to clarify the link between learning and transfer.

251 To summarize, we find an important role for perceptual regions in the prediction of future value-driven choice  
252 outcomes, which coincides with the sensitivity of BOLD in visual regions to differential value and signed  
253 feedback. These findings imply visual regions to learn prioritize high value features with the integration of  
254 feedback, to support and fasten, optimal response selection via the dorsal striatum in future encounters.

## 255 **Methods**

### 256 **Participants**

257 All participants had normal or corrected-to-normal vision and provided written consent before the scanning  
258 session, in accordance with the declaration of Helsinki. The ethics committee of the University of Amsterdam  
259 approved the experiment, and all procedures were in accordance with relevant laws and institutional guidelines.  
260 In total, six participants were excluded from all analyses due to movement (2), incomplete sessions (3), or  
261 misunderstanding of task instructions (1).

### 262 **Reinforcement learning task**

263 Full details of the reinforcement learning task are provided in 25. In brief, the task consisted of two phases  
264 (Figure 1a). In the first learning phase, three male or female face pairs (AB, CD, EF) were presented in a  
265 random order, and participants learned to select the most optimal face (A, C, E) in each pair solely through  
266 probabilistic feedback ('correct': happy smiley, 'incorrect': sad smiley). Choosing face-A lead to 'correct' on  
267 80% of the trials, whereas a choice for face-B only lead to the feedback 'correct' for 20% of the trials. Other  
268 ratios for 'correct' were 70:30 (CD) and 60:40 (EF). Participants were not informed about the complementary  
269 relationship in pairs. All trials started with a jitter interval where only a white fixation cross was presented  
270 and had a duration of 0, 500, 1000 or 1500ms to obtain an interpolated temporal resolution of 500ms. Two  
271 faces were then shown left and right of the fixation-cross and remained on screen up to response, or trial end  
272 (4000ms). If a response was given on time, a white box surrounding the chosen face was then shown (300ms)  
273 and followed (interval 0-450ms) by feedback (500ms). Omissions were followed by the text 'miss' (2000ms).  
274 The transfer-phase contained the three face-pairs from the learning phase, and 12 novel combinations, in which  
275 participants had to select which item they thought had been more rewarding during learning. Transfer-phase  
276 trials were identical to the learning phase, with the exception that no feedback was provided. All trials had  
277 a fixed duration of 4000ms, where in addition to the jitter used at the beginning of each trial, null trials  
278 (4000ms) were randomly interspersed across the learning (60 trials; 20%) and transfer (72 trials; 20%) phase.  
279 Each face was presented equally often on the left or right side, and choices were indicated with the right-hand  
280 index (left) or middle (right) finger. Before the MRI session, participants performed a complete learning  
281 phase to familiarize with the task (300 trials with different faces). In the MRI scanner, participants performed  
282 two learning blocks of 150 trials each (300 trials total; equal numbers of AB, CD and EF), and three transfer  
283 phase blocks of 120 trials each (360 total; 24 presentations of each pair). All stimuli were presented on a  
284 black-projection screen that was viewed via a mirror-system attached to the MRI head coil.

## 285 Reinforcement learning model

286 Trial-by-trial updating in value beliefs about the face selected in the learning phase, and reward predic-  
287 tion errors (signed expectancy violations) were estimated with a variant of the computational  $Q$ -learning  
288 algorithm<sup>27,64,65</sup> that is frequently used with this reinforcement learning task and contains two separate  
289 learning rate parameters for positive ( $\alpha_{gain}$ ) and negative ( $\alpha_{loss}$ ) reward prediction errors<sup>4,23,25,58</sup>.  $Q$ -learning  
290 assumes participants to maintain reward expectations for each of the six (A-to-F) stimuli presented during  
291 the learning phase. The expected value ( $Q$ ) for selecting a stimulus  $i$  (could be A-to-F) upon the next  
292 presentation is then updated as follows:

$$Q_i(t+1) = Q_i(t) + \begin{cases} \alpha_{Gain}[r_i(t) - Q_i(t)], & \text{if } r = 1 \\ \alpha_{Loss}[r_i(t) - Q_i(t)], & \text{if } r = 0 \end{cases}$$

293 Where  $0 \leq \alpha_{gain}$  or  $\alpha_{loss} \leq 1$  represent learning rates,  $t$  is trial number, and  $r = 1$  (positive feedback) or  $r = 0$   
294 (negative feedback). The probability of selecting one response over the other (i.e., A over B) is computed as:

$$P_A(t) = \frac{\exp(\beta * Q_t(A))}{\exp(\beta * Q_t(B)) + \exp(\beta * Q_t(A))}$$

295 With  $0 \leq \beta \leq 100$  known as the inverse temperature.

## 296 Bayesian hierarchical estimation procedure

297 To fit this  $Q$ -learning algorithm with two learning rate parameters we used Bayesian hierarchical estimation  
298 procedure. The full estimation procedure is explained in<sup>25</sup>. To summarize, this implementation assumes that  
299 probit-transformed model parameters for each participant are drawn from a group-level normal distribution  
300 characterized by group level mean and standard deviation parameters:  $z \sim N(\mu_z, \sigma_z)$ . A normal prior was  
301 assigned to group-level means  $\mu_z \sim N(0, 1)$ , and a uniform prior to the group-level standard deviations  
302  $\sigma_z \sim U(1, 1.5)$ . Model fits were implemented in Stan, where multiple chains were generated to ensure  
303 convergence.

## 304 Image acquisition

305 The fMRI data for the Reinforcement learning task was acquired in a single scanning session with two learning  
306 and three transfer phase runs on a 3-T scanner (Philips Achieva TX, Andover, MA) using a 32-channel head  
307 coil. Each scanning run contained 340 functional  $T2^*$ -weighted echo-planar images for the learning phase,  
308 and 290  $T2^*$ -weighted echo planar images for the transfer phase (TR = 2000 ms; TE = 27.63 ms; FA =  
309 76.1°; 3 mm slice thickness; 0.3 mm slice spacing; FOV = 240 × 121.8 × 240; 80 × 80 matrix; 37 slices,  
310 ascending slice order). After a short break of 10 minutes with no scanning, data collection was continued  
311 with a three-dimensional  $T1$  scan for registration purposes (repetition time [TR] = 8.5080 ms; echo time [TE]  
312 = 3.95ms; flip angle [FA] = 8°; 1 mm slice thickness; 0 mm slice spacing; field of view [FOV] = 240 × 220  
313 × 188), the fMRI data collection using a stop signal task (described in 25), and a localizer task with faces,  
314 houses, objects, and scrambled scenes to identify FFA responsive regions on an individual level (317  $T2^*$



315 weighted echo-planar images; TR = 1500 msec; TE = 27.6 msec; FA = 70°; 2.5 mm slice thickness; 0.25 mm  
316 slice spacing; FOV = 240 × 79.5 × 240; 96 × 96 matrix; 29 slices, ascending slice order). Here, participants  
317 viewed a series of houses, faces, objects as well as phase-scrambled scenes. To sustain attention during  
318 functional localization, subjects pressed a button when an image was directly repeated (12.5% likelihood).

## 319 **fMRI analysis learning phase**

320 The interplay between learning and perceptual activity was examined by evaluating how trial-by-trial  
321 computations of value-beliefs, and reward prediction errors relate to BOLD responses in the occipital cortex  
322 (OC) and fusiform face area (FFA). To compare perceptual responses with the more traditional literature, we  
323 first show how value-beliefs and RPEs relate to the activity pattern of the dorsal (i.e., caudate, or putamen)  
324 or ventral (i.e., accumbens) parts of the striatum. Regions of interest (ROI) templates were defined using  
325 anatomical atlases available in FSL, or the localizer task for FFA. For this purpose, the localizer scans were  
326 preprocessed using motion correction, slice-time correction, and pre-whitening<sup>66</sup>. For each subject, a GLM  
327 was fitted with the following EVs: for FFA, faces > (houses and objects), for PPA, houses > (faces and  
328 objects) and for LOC, intact scenes > scrambled scenes. Higher-level analysis was performed using FLAME  
329 Stage 1 and Stage 2 with automatic outlier detection<sup>67</sup>. For the whole-brain analysis Z (Gaussianized T/F)  
330 statistic images were thresholded using clusters determined by  $z > 2.3$  and  $p < .05$  (GRFT) to define a  
331 group-level binary FFA region. Templates used for the caudate [center of gravity (cog): (-) 13, 10, 10],  
332 putamen [cog: (-) 25, 1, 1], and Nucleus accumbens [cog: (-)19, 12, -7] were based on binary masks. Because  
333 participants were asked to differentiate faces, for each participant, we multiplied the binary templates of OC  
334 (V1) [cog: 1, -83, 5], FFA [cog: 23, -48, -18] with the individual t-stats from the localizer task contrast faces  
335 > (houses and objects). All anatomical masks, and the localizer group-level FFA mask can be downloaded  
336 from github (see acknowledgements).

## 337 **Deconvolution analysis learning phase**

338 To more precisely examine the time course of activation in the striatal and perceptual regions, we performed  
339 finite impulse response estimation (FIR) on the BOLD signals. After motion correction, temporal filtering  
340 (3rd order savitzky-golay filter with window of 120 s) and percent signal change conversion, data from each  
341 region was averaged across voxels while weighting voxels according to ROI probability masks, and upsampled  
342 from 0.5 to 3 Hz. This allows the FIR fitting procedure to capitalize on the random timings (relative to  
343 TR onset) of the stimulus presentation and feedback events in the experiment. Separate response time  
344 courses were simultaneously estimated triggered on two separate events: stimulus onset, feedback onset. FIR  
345 time courses for all trial types were estimated simultaneously using a penalized (ridge) least-squares fit, as  
346 implemented in the FIRDeconvolution package<sup>68</sup>, and the appropriate penalization parameter was estimated  
347 using cross-validation. For stimulus onset events (i.e., onset presentation of face pairs) response time courses  
348 were fit separately for the AB, CD and EF pairs, while also estimating the time courses of signal covariation  
349 with chosen and unchosen value for these pairs. For these events, our analysis corrected for the duration of  
350 the decision process. For the feedback events, the co-variation response time course with signed and unsigned  
351 prediction errors were estimated. These signal response time courses were analysed using across-subjects  
352 GLMs at each time-point using the statsmodels package<sup>69</sup>. The  $\alpha$  value for the contributions of  $Q$  or RPE  
353 was set to 0.0125 (i.e. a Bonferroni corrected value of 0.05 given the interval of interest between 0 and 8 s).



## 354 **Random Forest classification**

355 To specify the relevance of perceptual regions in the resolve of future value-driven choices a random forest  
356 (RF) classifier was used<sup>28</sup>. The RF classifier relies on an ensemble of decision trees as base learners, where  
357 the final prediction (e.g., for a given trial is the choice going to be correct/optimal? or incorrect/suboptimal?  
358 given past learning) is obtained by a majority vote that combines the prediction of all decision trees. To  
359 achieve controlled variation, each decision tree is trained on a random subset of the variables (i.e. regions  
360 of interest chosen), and a bootstrapped sample of data points (i.e. trials). In the construction of each tree  
361 about 1/3 of all trials is left out - termed as the out-of-bag sample – and later used to see how well each tree  
362 preforms on unseen data. The generalized error for predictions is calculated by aggregating the prediction for  
363 every out-of-bag sample across all trees. An important feature of the RF classification method is the ease to  
364 measure the relative importance of each variable (i.e., region), in the overall predictive performance. That is,  
365 it allows for the ranking of all regions evaluated in the prediction of future value-based decisions.

## 366 **ROI selection and Random Forest procedure**

367 This study used the ‘Breiman and Cutler’s Random Forests for Classification and Regression’ package in R,  
368 termed randomForest. RF evaluations relied on the fMRI data recorded during the transfer phase, in a set  
369 of 9 regions of interest (ROIs). These ROIs included all templates from the learning phase (i.e., caudate,  
370 putamen, accumbens, OC, and FFA), as well as, the ventromedial prefrontal cortex (vmPFC), dorsolateral  
371 prefrontal cortex (DLPFC), pre-supplementary motor area (preSMA), and the primary motor cortex (M1).  
372 The selection of these additional anatomical templates was inspired by our previous analysis of this data  
373 with those templates focusing on networks<sup>25,62,70</sup>. From each ROI a single parameter estimate (averaged  
374 normalized  $\beta$  estimate across voxels in each ROI) was obtained per trial, per subject. All, pre-processing steps  
375 to obtain single-trial images are described in 25. Single-trial activity estimates were used as input variables  
376 in RF to predict choice outcomes (optimal/sub-optimal) in the transfer phase. Here, participants choose the  
377 best/optimal option based on values learned during the learning phase. We defined optimal choices as correct  
378 (i.e, when participants choose the option with the higher value), and sub-optimal choices as incorrect. Misses  
379 were excluded from RF evaluations.

380 By design, the transfer-phase contained 360 trials including 15 different pairs (12 novel), where each pair was  
381 presented 24 times with the higher value presented left in 12 of the 24 presentations, and on the right for the  
382 other half. With so many subtle value differences across the options presented and only one BOLD estimate  
383 per trial/region the prediction of future choices is under powered (Figure 5a). Therefore, assuming that all  
384 participants come from the same population, a fixed effects approach was taken for evaluations with RF. Here,  
385 the trial\*region activity matrices for all participants were combined into one big data matrix (Figure 5b)  
386 and subsequently shuffled across the rows, so that both participants and trials were re-arranged in a random  
387 order across rows. Besides the single trial BOLD estimates from the 9 ROI’s, this shuffled matrix contained  
388 two additional columns, which specified subject\_id (to which subject does each trial belong), and Trial Sign –  
389 i.e., is the choice between the two faces about two positive (+/+; AC, AE, CE), negative (-/-; BD, BF, DF),  
390 or a positive-negative (+/-; e.g. AD, CF etc. ) associations given the task manipulation during learning.  
391 Subject\_id was included to control for different BOLD fluctuations across participants, whereas Trial Sign  
392 was added because both BOLD and choice patterns differ across these options (please see 25). The shuffled  
393 fixed effect matrix was divided into a separate training (2/3 of whole matrix), and validation (1/3) set, to

394 be used for RF evaluations (Figure 5c). Learning was based on the training set, using 2000 trees with the  
395 number of variables (regions) used by each tree optimized with the tuneRF function in R, and accordingly  
396 set to 5. For the construction of each tree about 1/3 of all trials is left out - termed as the out-of-bag sample  
397 - and later used to see how well each tree performs on unseen data. The generalized error for predictions is  
398 calculated by aggregating the prediction for every out-of-bag sample across all trees. Besides this out-of-bag  
399 approximation we evaluated the predictive accuracy of the whole RF on the separate unseen validation-set.  
400 The single trial data used as input, the RF evaluation codes, and ROI templates can all be downloaded from  
401 the github link provided in acknowledgements.

## 402 **Acknowledgements**

403 This work was supported by an ABC Talent grant to SJ from the University of Amsterdam, an ERC grant  
404 ERC-2012-AdG-323413 to JT, and NWO-CAS grant 012.200.012 to TK.

## 405 **Author contribution**

406 SJ and TK developed the questions and analysis plan for the re-analysis. SJ and TK contributed novel  
407 methods and analyzed the data. SJ wrote the first draft of the MS with edits from TK. JT commented on  
408 the final draft.

## 409 **Competing interests**

410 The authors declare to have no competing interests.

## 411 **Data availability**

412 The code and preprocessed files for behavioral and decoding analyses can be download from: [https://](https://github.com/sarajahfari/Pearl3T)  
413 [github.com/sarajahfari/Pearl3T](https://github.com/sarajahfari/Pearl3T).git, and fMRI preprocessing and deconvolution analysis code are available at  
414 [https://github.com/tnapen/pearl\\_3T](https://github.com/tnapen/pearl_3T). The raw data can be downloaded from openfMRI in BIDS after  
415 acceptance of this MS.

## 416 References

- 417 1. Daw, N. D., O’Doherty, J. P., Dayan, P., Seymour, B. & Dolan, R. J. Cortical substrates for exploratory  
418 decisions in humans. *Nature* **441**, 876–879 (2006).
- 419 2. Hare, T. A., Schultz, W., Camerer, C. F., O’Doherty, J. P. & Rangel, A. Transformation of stimulus value  
420 signals into motor commands during simple choice. *Proceedings of the National Academy of Sciences* **108**,  
421 18120–18125 (2011).
- 422 3. Jocham, G., Klein, T. A. & Ullsperger, M. Dopamine-mediated reinforcement learning signals in the  
423 striatum and ventromedial prefrontal cortex underlie value-based choices. *Journal of Neuroscience* **31**,  
424 1606–1613 (2011).
- 425 4. Kahnt, T. *et al.* Dorsal striatal–midbrain connectivity in humans predicts how reinforcements are used to  
426 guide decisions. *Journal of Cognitive Neuroscience* **21**, 1332–1345 (2009).
- 427 5. Klein, T. A., Ullsperger, M. & Jocham, G. Learning relative values in the striatum induces violations of  
428 normative decision making. *Nature Communications* **8**, 16033 (2017).
- 429 6. O’Doherty, J. P., Cockburn, J. & Pauli, W. M. Learning, reward, and decision making. *Annual review of*  
430 *psychology* **68**, 73–100 (2017).
- 431 7. O’Doherty, J. *et al.* Dissociable roles of ventral and dorsal striatum in instrumental conditioning. *Science*  
432 **304**, 452–454 (2004).
- 433 8. Schultz, W., Dayan, P. & Montague, P. R. A neural substrate of prediction and reward. *Science* **275**,  
434 1593–1599 (1997).
- 435 9. Montague, P. R., Dayan, P. & Sejnowski, T. J. A framework for mesencephalic dopamine systems based  
436 on predictive hebbian learning. *Journal of neuroscience* **16**, 1936–1947 (1996).
- 437 10. Tobler, P. N., Fiorillo, C. D. & Schultz, W. Adaptive coding of reward value by dopamine neurons.  
438 *Science* **307**, 1642–1645 (2005).
- 439 11. Atallah, H. E., Lopez-Paniagua, D., Rudy, J. W. & O’Reilly, R. C. Separate neural substrates for skill  
440 learning and performance in the ventral and dorsal striatum. *Nature neuroscience* **10**, 126–131 (2007).
- 441 12. Joel, D., Niv, Y. & Ruppin, E. Actor–critic models of the basal ganglia: New anatomical and computational  
442 perspectives. *Neural networks* **15**, 535–547 (2002).
- 443 13. Collins, A. G. E. & Frank, M. J. Opponent actor learning (opal): Modeling interactive effects of striatal  
444 dopamine on reinforcement learning and choice incentive. *Psychological review* **121**, 337–366 (2014).
- 445 14. Hikosaka, O., Kim, H. F., Yasuda, M. & Yamamoto, S. Basal ganglia circuits for reward value–guided  
446 behavior. *Annual review of neuroscience* **37**, 289–306 (2014).
- 447 15. Kim, H. F. & Hikosaka, O. Distinct basal ganglia circuits controlling behaviors guided by flexible and  
448 stable values. *Neuron* **79**, 1001–1010 (2013).
- 449 16. Kim, H. F., Amita, H. & Hikosaka, O. Indirect pathway of caudal basal ganglia for rejection of valueless  
450 visual objects. *Neuron* **94**, 920–930 (2017).
- 451 17. Serences, J. T. Value-based modulations in human visual cortex. *Neuron* **60**, 1169–1181 (2008).

- 452 18. Serences, J. T. & Saproo, S. Population response profiles in early visual cortex are biased in favor of  
453 more valuable stimuli. *Journal of neurophysiology* **104**, 76–87 (2010).
- 454 19. Kravitz, D. J., Saleem, K. S., Baker, C. I., Ungerleider, L. G. & Mishkin, M. The ventral visual pathway:  
455 An expanded neural framework for the processing of object quality. *Trends in cognitive sciences* **17**, 26–49  
456 (2013).
- 457 20. Fernandez-Ruiz, J., Wang, J., Aigner, T. G. & Mishkin, M. Visual habit formation in monkeys with  
458 neurotoxic lesions of the ventrocaudal neostriatum. *Proceedings of the National Academy of Sciences* **98**,  
459 4196–4201 (2001).
- 460 21. Lim, S.-L., O’Doherty, J. P. & Rangel, A. The decision value computations in the vmPFC and striatum  
461 use a relative value code that is guided by visual attention. *Journal of Neuroscience* **31**, 13214–13223 (2011).
- 462 22. Lim, S.-L., O’Doherty, J. P. & Rangel, A. Stimulus value signals in ventromedial pfc reflect the integration  
463 of attribute value signals computed in fusiform gyrus and posterior superior temporal gyrus. *Journal of*  
464 *Neuroscience* **33**, 8729–8741 (2013).
- 465 23. Jahfari, S. & Theeuwes, J. Sensitivity to value-driven attention is predicted by how we learn from value.  
466 *Psychonomic bulletin & review* **24**, 408–415 (2017).
- 467 24. Jahfari, S., Waldorp, L., Ridderinkhof, K. R. & Scholte, H. S. Visual information shapes the dynamics of  
468 corticobasal ganglia pathways during response selection and inhibition. *Journal of cognitive neuroscience* **27**,  
469 1344–1359 (2015).
- 470 25. Jahfari, S. *et al.* Cross-task contributions of frontobasal ganglia circuitry in response inhibition and  
471 conflict-induced slowing. *Cerebral Cortex* bhy076 (2018).
- 472 26. O’Doherty, J. P., Hampton, A. & Kim, H. Model-based fMRI and its application to reward learning and  
473 decision making. *Annals of the New York Academy of sciences* **1104**, 35–53 (2007).
- 474 27. Daw, N. D. Trial-by-trial data analysis using computational models. *Decision making, affect, and learning:*  
475 *Attention and performance XXIII* **23**, 3–38 (2011).
- 476 28. Breiman, L. Random forests. *Machine learning* **45**, 5–32 (2001).
- 477 29. Breiman, L. Consistency for a simple model of random forests. (2004).
- 478 30. Fouragnan, E., Retzler, C. & Philiastides, M. G. Separate neural representations of prediction error  
479 valence and surprise: Evidence from an fMRI meta-analysis. *Human brain mapping* (2018).
- 480 31. Cicmil, N., Cumming, B. G., Parker, A. J. & Krug, K. Reward modulates the effect of visual cortical  
481 microstimulation on perceptual decisions. *Elife* **4**, e07832 (2015).
- 482 32. Shuler, M. G. & Bear, M. F. Reward timing in the primary visual cortex. *Science* **311**, 1606–1609 (2006).
- 483 33. Pleger, B. *et al.* Influence of dopaminergically mediated reward on somatosensory decision-making. *PLoS*  
484 *biology* **7**, e1000164 (2009).
- 485 34. Kaskan, P. M. *et al.* Learned value shapes responses to objects in frontal and ventral stream networks in  
486 macaque monkeys. *Cerebral Cortex* **27**, 2739–2757 (2016).

- 487 35. Kahnt, T., Grueschow, M., Speck, O. & Haynes, J.-D. Perceptual learning and decision-making in human  
488 medial frontal cortex. *Neuron* **70**, 549–559 (2011).
- 489 36. Den Ouden, H. E. M., Kok, P. & De Lange, F. P. How prediction errors shape perception, attention, and  
490 motivation. *Frontiers in psychology* **3**, 548 (2012).
- 491 37. Gottlieb, J. Attention, learning, and the value of information. *Neuron* **76**, 281–295 (2012).
- 492 38. Gottlieb, J., Hayhoe, M., Hikosaka, O. & Rangel, A. Attention, reward, and information seeking. *Journal*  
493 *of Neuroscience* **34**, 15497–15504 (2014).
- 494 39. Störmer, V., Eppinger, B. & Li, S.-C. Reward speeds up and increases consistency of visual selective  
495 attention: A lifespan comparison. *Cognitive, Affective, & Behavioral Neuroscience* **14**, 659–671 (2014).
- 496 40. Van Slooten, J. C., Jahfari, S., Knapen, T. & Theeuwes, J. How pupil responses track value-based  
497 decision-making during and after reinforcement learning. *PLoS computational biology* **14**, e1006632 (2018).
- 498 41. McCoy, B., Jahfari, S., Engels, G., Knapen, T. & Theeuwes, J. Dopaminergic medication reduces striatal  
499 sensitivity to negative outcomes in parkinson’s disease. *bioRxiv* (2018).
- 500 42. Lak, A., Stauffer, W. R. & Schultz, W. Dopamine neurons learn relative chosen value from probabilistic  
501 rewards. *Elife* **5**, e18044 (2016).
- 502 43. Lak, A., Nomoto, K., Keramati, M., Sakagami, M. & Kepecs, A. Midbrain dopamine neurons signal belief  
503 in choice accuracy during a perceptual decision. *Current Biology* **27**, 821–832 (2017).
- 504 44. Ding, L. & Gold, J. I. Caudate encodes multiple computations for perceptual decisions. *Journal of*  
505 *Neuroscience* **30**, 15747–15759 (2010).
- 506 45. Yamamoto, S., Monosov, I. E., Yasuda, M. & Hikosaka, O. What and where information in the caudate  
507 tail guides saccades to visual objects. *Journal of Neuroscience* **32**, 11005–11016 (2012).
- 508 46. Hikosaka, O., Yamamoto, S., Yasuda, M. & Kim, H. F. Why skill matters. *Trends in cognitive sciences*  
509 **17**, 434–441 (2013).
- 510 47. Hassabis, D., Kumaran, D., Summerfield, C. & Botvinick, M. Neuroscience-inspired artificial intelligence.  
511 *Neuron* **95**, 245–258 (2017).
- 512 48. Hebart, M. N. & Baker, C. I. Deconstructing multivariate decoding for the study of brain function.  
513 *Neuroimage* **180**, 4–18 (2018).
- 514 49. Naselaris, T., Kay, K. N., Nishimoto, S. & Gallant, J. L. Encoding and decoding in fMRI. *Neuroimage*  
515 **56**, 400–410 (2011).
- 516 50. Snoek, L., Miletic, S. & Scholte, H. S. How to control for confounds in decoding analyses of neuroimaging  
517 data. *NeuroImage* **184**, 741–760 (2019).
- 518 51. Kriegeskorte, N. & Douglas, P. K. Interpreting encoding and decoding models. *arXiv preprint*  
519 *arXiv:1812.00278* (2018).
- 520 52. King, J.-R. *et al.* Encoding and decoding neuronal dynamics: Methodological framework to uncover the  
521 algorithms of cognition. (2018).

- 522 53. Cools, R. & D'Esposito, M. Inverted-u-shaped dopamine actions on human working memory and cognitive  
523 control. *Biological psychiatry* **69**, e113–e125 (2011).
- 524 54. Aston-Jones, G. & Cohen, J. D. An integrative theory of locus coeruleus-norepinephrine function:  
525 Adaptive gain and optimal performance. *Annu. Rev. Neurosci.* **28**, 403–450 (2005).
- 526 55. Yu, A. J. & Dayan, P. Uncertainty, neuromodulation, and attention. *Neuron* **46**, 681–692 (2005).
- 527 56. Beste, C. *et al.* Dopamine modulates the efficiency of sensory evidence accumulation during perceptual  
528 decision making. *International Journal of Neuropsychopharmacology* (2018).
- 529 57. O'Doherty, J., Critchley, H., Deichmann, R. & Dolan, R. J. Dissociating valence of outcome from  
530 behavioral control in human orbital and ventral prefrontal cortices. *Journal of neuroscience* **23**, 7931–7939  
531 (2003).
- 532 58. Niv, Y., Edlund, J. A., Dayan, P. & O'Doherty, J. P. Neural prediction errors reveal a risk-sensitive  
533 reinforcement-learning process in the human brain. *Journal of Neuroscience* **32**, 551–562 (2012).
- 534 59. Sasikumar, D., Emeric, E., Stuphorn, V. & Connor, C. E. First-pass processing of value cues in the  
535 ventral visual pathway. *Current Biology* **28**, 538–548 (2018).
- 536 60. Jocham, G., Boorman, E. & Behrens, T. Neuroscience of value-guided choice. *The Wiley Handbook on*  
537 *the Cognitive Neuroscience of Learning* 554–591 (2016).
- 538 61. Shenhav, A., Straccia, M. A., Cohen, J. D. & Botvinick, M. M. Anterior cingulate engagement in a  
539 foraging context reflects choice difficulty, not foraging value. *Nature neuroscience* **17**, 1249 (2014).
- 540 62. Pircalabelu, E., Claeskens, G., Jahfari, S. & Waldorp, L. J. A focused information criterion for graphical  
541 models in fMRI connectivity with high-dimensional data. *The Annals of Applied Statistics* **9**, 2179–2214  
542 (2015).
- 543 63. Bhandari, A., Gagne, C. & Badre, D. Just above chance: Is it harder to decode information from human  
544 prefrontal cortex blood oxygenation level-dependent signals? *Journal of cognitive neuroscience* 1–26 (2018).
- 545 64. Frank, M. J., Moustafa, A. A., Haughey, H. M., Curran, T. & Hutchison, K. E. Genetic triple dissociation  
546 reveals multiple roles for dopamine in reinforcement learning. *Proceedings of the National Academy of Sciences*  
547 **104**, 16311–16316 (2007).
- 548 65. Watkins, C. J. C. H. & Dayan, P. Q-learning. *Machine learning* **8**, 279–292 (1992).
- 549 66. Woolrich, M. W., Ripley, B. D., Brady, M. & Smith, S. M. Temporal autocorrelation in univariate linear  
550 modeling of fMRI data. *Neuroimage* **14**, 1370–1386 (2001).
- 551 67. Beckmann, C. F., Jenkinson, M. & Smith, S. M. General multilevel linear modeling for group analysis in  
552 fmri. *Neuroimage* **20**, 1052–1063 (2003).
- 553 68. Knapen, T. & JW, G. FIRDeconvolution. (2016). doi:doi:10.5281/ZENODO.46216
- 554 69. Seabold, S. & Perktold, J. Statsmodels: Econometric and statistical modeling with python. in *Proceedings*  
555 *of the 9th python in science conference* **57**, 57–61 (2010).
- 556 70. Schmittmann, V. D., Jahfari, S., Borsboom, D., Savi, A. O. & Waldorp, L. J. Making large-scale networks  
557 from fMRI data. *PloS one* **10**, e0129074 (2015).