1 # Title: Biased credit assignment in
2 # motivational learning biases arises
3 # through prefrontal influences on striatal
4 # learning

5

6 Short title: Motivational biases in fronto-striatal circuits

7 ## Authors

8

9 Johannes Algermissen[1]*, Jennifer C. Swart[1], René Scheeringa[1,2], Roshan Cools[1,3], Hanneke E.M. den
10 Ouden[1]*

11 ## Affiliations
12 [1] Radboud University, Donders Institute for Brain, Cognition and Behaviour, Nijmegen, The
13 Netherlands
14 [2] Erwin L. Hahn Institute for Magnetic Resonance Imaging, University of Duisburg-Essen, Essen,
15 Germany
16 [3] Department of Psychiatry, Radboud University Medical Centre, Nijmegen, The Netherlands
17 * j.algermissen@donders.ru.nl; h.denouden@donders.ru.nl

18 ## Abstract
19 Actions are biased by the outcomes they can produce: Humans are more likely to show action under
20 reward prospect, but hold back under punishment prospect. Such motivational biases derive not only
21 from biased response selection, but also from biased learning: humans tend to attribute rewards to
22 their own actions, but are reluctant to attribute punishments to having held back. The neural origin
23 of these biases is unclear; in particular, it remains open whether motivational biases arise solely from
24 an evolutionarily old, subcortical architecture or also due to younger, cortical influences.
25 Simultaneous EEG-fMRI allowed us to track which regions encoded biased prediction errors in which
26 order. Biased prediction errors occurred in cortical regions (ACC, vmPFC, PCC) before subcortical
27 regions (striatum). These results highlight that biased learning is not a mere feature of the basal
28 ganglia, but arises through prefrontal cortical contributions, revealing motivational biases to be a
29 potentially flexible, sophisticated mechanism.

30 ## Teaser

31

32 Cortical influences on subcortical learning explain why we attribute rewards to actions, but not
33 punishments to inactions.

## Introduction

Human action selection is biased by potential action outcomes: reward prospect drives us to invigorate action, while threat of punishment holds us back (*1–3*). These motivational biases have been evoked to explain why humans are tempted by reward-related cues signaling the chance to gain food, drugs, or money, as they elicit automatic approach behavior. Conversely, punishment-related cues suppress action and lead to paralysis, which may even lie at the core of mental health problems such as phobias and mood disorders (*4, 5*). While such examples highlight the potential maladaptiveness of biases in some situations, they confer benefits in other situations: Biases could provide sensible "default" actions before context-specific knowledge is acquired (*1, 6*). They may also provide ready-made alternatives to more demanding action selection mechanisms, especially when speed has to be prioritized (*7*).

Previous research has assumed that motivational biases arise because the valence of prospective outcomes influences action selection (*8*). However, we have recently shown that not only action selection, but also the updating of action values based on obtained outcomes is subject to valence-dependent biases (*3, 9, 10*): humans are more inclined to ascribe rewards to active responses, but have problems with attributing punishments to having held back. One the one hand, such biased learning might be adaptive in combining the flexibility of instrumental learning with somewhat rigid "priors" about typical action-outcome relationships. Exploiting lifetime (or evolutionary) experience might lead to learning that is faster and more robust to environmental "noise". On the other hand, biases might be responsible for phenomena of "animal superstition" like auto-shaping or negative maintenance, where rats and pigeons repeat behavioral patterns that co-occurred with the attainment of (factually random) rewards and keep showing such behavior even if it delays or decreases rewards (*1, 11, 12*). While reward attainment can lead to an illusory sense of control over outcomes, control is underestimated under threat of punishment: Humans find it hard to comprehend how inactions can cause negative outcomes, which makes them more lenient in judging harms caused by others' inactions (*13, 14*). Taken together, also credit assignment is subject to motivational biases, with enhanced credit for rewards given to actions, but diminished credit for punishments given to inactions.

While evident in behavior, the neural mechanisms subserving such biased credit assignment are unclear. One strong candidate region is the striatum, part of the evolutionarily old basal ganglia system. Influential computational models of basal ganglia function (*15, 16*) (henceforth called "asymmetric pathways model") predict such motivational learning biases: Positive prediction errors, elicited by rewards, lead to long-term potentiation in the striatal direct "Go" pathway (and long term depression in the indirect pathway), allowing for a particularly effective acquisition of Go responses after rewards. Conversely, negative prediction errors, elicited by punishments, lead to long term potentiation in the "NoGo" pathway, impairing the unlearning of NoGo responses after punishments. This account suggests that motivational biases arise within the same pathways involved in standard reinforcement learning (RL). An alternative candidate model is that biases arise through the modulation of these evolutionarily old RL systems by external, evolutionarily younger areas that also track past actions, putatively the prefrontal cortex (PFC). Past research has suggested that standard RL can be biased by information stored in PFC, such as explicit instructions (*17, 18*) or cognitive-map like models of the environment (*19–21*). Most notably, the anterior cingulate cortex (ACC) has been found to reflect the impact of explicit instructions (*18*) and of environmental changes on prediction errors (*22, 23*).

78      Both candidate models predict that BOLD signal in striatum should be better described by
79  biased compared with "standard" prediction errors. In addition, the model proposing a prefrontal
80  influence on striatal processing makes a notable prediction about the timing of signals: information
81  about the selected action and the obtained outcome should be present first in prefrontal circuits to
82  then later affect processes in the striatum. While fMRI BOLD recordings allow for unequivocal access
83  to striatal activity, the sluggish nature of the BOLD signal prevents clear inferences about temporal
84  precedence of signals from different regions. We thus combined BOLD with simultaneous EEG
85  recordings which allowed us to precisely characterize learning signals in both space and time.
86      The key question is whether biased credit assignment arises directly from biased RL through
87  the asymmetric pathways in the striatum, or whether striatal RL mechanisms are biased by external
88  prefrontal sources, with the ACC as likely candidate. To this end, participants performed a
89  motivational Go/ NoGo learning task that is well-established to evoke motivational biases of action
90  (*3, 9, 24*). We expected to observe biased PEs in striatum and frontal cortical areas. By
91  simultaneously recording fMRI and EEG and correlating trial-by-trial BOLD signal with EEG time-
92  frequency power, we were able to time-lock the peaks of EEG-BOLD correlations for regions
93  reflecting biased PEs and infer their relative temporal precedence. We focused on two well-
94  established electrophysiological signatures of RL, namely theta and delta power (*25–30*) as well as
95  beta power (*25, 31*) over midfrontal electrodes.

## Results

97      Thirty-six participants performed a motivational Go/ NoGo learning task (*3, 9*) in which required
98  action (Go/ NoGo) and potential outcome (reward/ punishment) were orthogonalized (Fig. 1A-D).
99  They learned by trial-and-error for each of eight cues whether to perform a left button press ($Go_{LEFT}$),
100  right button press ($Go_{RIGHT}$), or no button press (NoGo), and whether a correct action increased the
101  chance to win a reward (Win cues) or to avoid a punishment (Avoid cues). Correct actions lead to
102  80% favorable outcomes (reward, no punishment), with only 20% favorable outcomes for incorrect
103  actions. Participants performed two sessions of 320 trials, with separate cue sets, which were
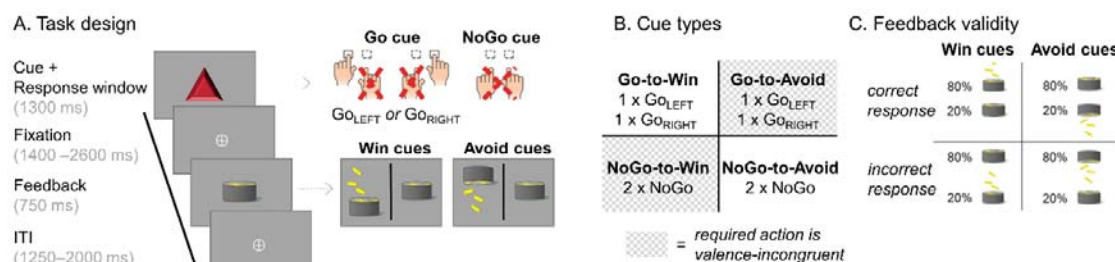104  counterbalanced across participants.
105



**Figure 1. Motivational Go/ NoGo learning task design.** (A) On each trial, a Win or Avoid cue appears; valence of the cue is not signaled but should be learned. Cue offset is also the response deadline. Response-dependent feedback follows after a jittered interval. Each cue has only one correct action ($Go_{LEFT}$, $Go_{Right}$, or NoGo), which is followed by the favorable outcome 80% of the time. For Win cues, actions can lead to rewards or neutral outcomes; for Avoid cues, actions can lead to neutral outcomes or punishment. Rewards and punishments are depicted by money falling into/ out of a can. (B) There are eight different cues, orthogonalizing cue valence (Win versus Avoid) and required action (Go versus NoGo). The motivationally incongruent cues, for which the motivational action tendencies are incongruent with the instrumental requirements, are highlighted in gray. (C) Feedback is probabilistic: Correct actions to Win cues lead to rewards in 80% of cases, but neutral outcomes in 20% of cases. For Avoid cues, correct actions lead to neutral outcomes in 80% of cases, but punishments in

20% of cases. For incorrect actions, these probabilities are reversed.

106

## Regression analyses of behavior

107

108      We performed regression analyze to test whether a) responses were biased by the valence of
109      prospective outcomes (Win/ Avoid), reflecting biased responding and/ or learning, and b) whether
110      response repetition after favorable vs. non-favorable outcomes was biased by whether a Go vs.
111      NoGo response was performed, selectively reflecting biased learning.

112      For the first purpose, we analyzed choice data (Go/ NoGo) using mixed-effects logistic
113      regression that included factors required action (Go/ NoGo; note that this approach collapses across
114      $Go_{LEFT}$ and $Go_{RIGHT}$ responses), cue valence (Win/ Avoid), and their interaction (also reported in)(32).
115      Participants learned the task, i.e., they performed more Go responses towards Go than NoGo cues
116      (main effect of required action: $b = 0.815$, $SE = 0.113$, $\chi^2(1) = 32.008$, $p < .001$). In contrast to previous
117      studies (3, 9), learning did not asymptote (Fig. 2A), which provided greater dynamic range for the
118      biased learning effects to surface. Furthermore, participants showed a motivational bias, i.e., they
119      performed more Go responses to Win than Avoid cues (main effect of cue valence, $b = 0.423$, $SE =$
120      $0.073$, $\chi^2(1) = 23.695$, $p < .001$). Replicating other studies with this task, there was no significant
121      interaction between required action and cue valence ($b = 0.030$, $SE = 0.068$, $\chi^2(1) = 0.196$, $p = .658$,
122      Fig. 2A-B), i.e., there was no evidence for the effect of cue valence (motivational bias) differing in size
123      between Go or NoGo cues.

124      Secondly, as a proxy of (biased) learning, we analyzed cue-based response repetition
125      (probability of repeating a response on the next encounter of the same cue) as a function of outcome
126      valence (favorable vs non-favorable outcome), performed action (Go vs. NoGo), and outcome
127      salience (salient: reward or punishment vs. neutral: no reward or no punishment). As expected,
128      people were more likely to repeat the same response following a favorable outcome (main effect of
129      outcome valence: $b = 0.504$, $SE = 0.053$, $\chi^2(1) = 45.595$, $p < .001$). Most importantly, after salient
130      outcomes, participants adjusted their responses to a larger degree following Go responses than
131      NoGo responses, revealing the presence of a learning bias (Fig. 2C; interaction of valence x action x
132      salience: $b = 0.248$, $SE = 0.048$, $\chi^2(1) = 19.732$, $p < .001$). When selectively analyzing trials with salient
133      outcomes only, rewards (compared to punishments) led to a higher proportion of choice repetitions
134      following Go relative to NoGo responses (valence x response: $b = 0.308$, $SE = 0.064$, $\chi^2(1) = 17.798$, $p$
135      $< .001$; valence effect for Go only: $b = 1.276$, $SE = 0.115$, $\chi^2(1) = 53.932$, $p < .001$; valence effect for
136      NoGo only: $b = 0.637$, $SE = 0.127$, $\chi^2(1) = 18.228$, $p < .001$; see full results in S02).

137      Taken together, these results suggest that behavioral adaptation following rewards and
138      punishments is biased by the type of action that led to this outcome (Go or NoGo). However, these
139      analyses only consider behavioral adaptation on the next trial, and cannot pinpoint the precise
140      algorithmic nature of this learning bias. More importantly, it does not provide trial-by-trial estimates
141      of action values as required for model-based fMRI and EEG analyses to test for regions or time points
142      that reflect biased learning. We thus analyzed the impact of past outcomes on participants' choices
143      using computational RL models.

## Computational modeling of behavior

144

145      In line with previous work (3, 9), we fitted a series of increasingly complex RL models. We started
146      with a simple Rescorla Wagner model featuring learning rate and feedback sensitivity parameters
147      (M1). We next added a Go bias, capturing participants' overall propensity to make Go responses
148      (M2), and a Pavlovian response bias (M3), reflecting participants' propensity to adjust their likelihood

149   of emitting a Go response in response to Win vs. Avoid cues (*3*). Alternatively, we added an
150   instrumental learning bias (M4), amplifying the learning rate after rewarded Go responses and
151   dampening it after punished NoGo responses (*3*), in line with the asymmetric pathways model. In the
152   final model (M5), we added both a response bias and a learning bias. For the full model space (M1-
153   M5) and model definitions, see the Methods section. For a comparison with an alternative learning
154   bias specification based on the idea that active responses enhance credit assignment (*33*), see S04.
155       Model comparison showed clear evidence in favor of the full asymmetric pathways model
156   featuring both response and learning biases (M5; model frequency: 86.43%, protected exceedance
157   probability: 100%, see Fig. 2D, H; for model parameters and fit indices, see S03). Posterior predictive
158   checks involving one-step-ahead predictions and model simulations showed that this model captured
159   key behavioral features (Fig. 2E, F), including motivational biases and a greater behavioral adaptation
160   after Go responses followed by salient outcomes than after NoGo responses followed by salient
161   outcomes (Fig. 2 G). This pattern could not be captured by the alternative learning bias model (S04).
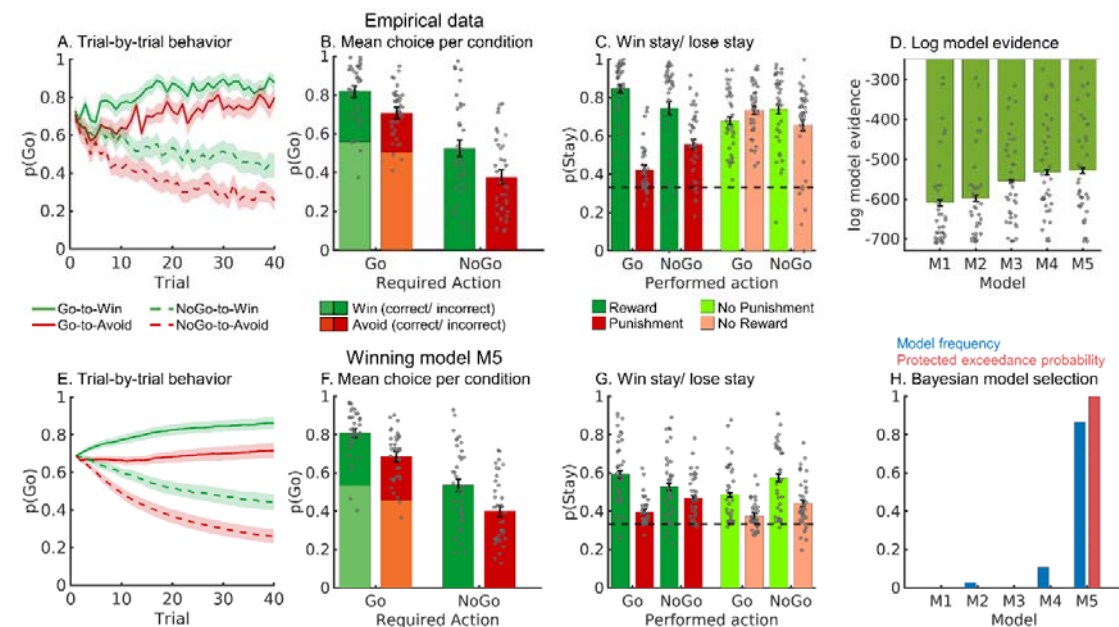162



**Figure 2. Behavioral performance.** (A) Trial-by-trial proportion of Go responses (±SEM across participants) for Go cues (solid lines) and NoGo cues (dashed lines). The motivational bias is already present from very early trials onwards, as participants made more Go responses to Win than Avoid cues (i.e., green lines are above red lines). Additionally, participants clearly learn whether to make a Go response or not (proportion of Go responses increases for Go cues and decreases for NoGo cues). (B) Mean (±SEM across participants) proportion Go responses per cue condition (points are individual participants' means). (C) Probability to repeat a response ("stay") on the next encounter of the same cue as a function of action and outcome. Learning is reflected in higher probability of staying after positive outcomes than after negative outcomes (main effect of outcome valence). Biased learning is evident in learning from salient outcomes, where this valence effect was stronger after Go responses than NoGo responses. Dashed line indicates chance level choice ($p_{stay}$ = 0.33). (D) Log-model evidence favors the asymmetric pathways model (M5) over simpler models (M1-M4). (E-G) Trial-by-trial proportion of Go responses, mean proportion Go responses, and probability of staying based on one-step-ahead predictions using parameters (hierarchical Bayesian inference) of the winning model (asymmetric pathways model, M5). (H) Model frequency and protected exceedance probability indicate best fit for model M5 (asymmetric pathways model), in line with log model evidence.

163

## fMRI: Basic quality control analyses

164

165    First, we performed a GLM as a quality-check to test which regions encoded favorable
166    (rewards, no punishments) vs. unfavorable (no reward/ punishment) outcomes in a "model-free"
167    way, independent of any model-based measure derived from a RL model (for full description of the
168    GLM regressors and contrasts, see S06). Favorable outcomes elicited a higher BOLD response in
169    regions including ventromedial PFC (vmPFC), ventral striatum, and right hippocampus, while
170    unfavorable outcomes elicited higher BOLD in bilateral dorsolateral PFC (dlPFC), left ventrolateral
171    PFC, and precuneous (Fig. 3A, see full report of significant clusters in S07).
172    We also assessed which regions encoded Go vs. NoGo as well Go$_{LEFT}$ vs. Go$_{RIGHT}$ responses.
173    There was higher BOLD for Go than NoGo responses at the time of response in PFC, ACC, striatum,
174    thalamus, motor cortices, and cerebellum, while BOLD was higher for NoGo than Go responses in
175    right IFG (Fig. 6C left panel; see S04)(32). For lateralized Go responses, there was higher BOLD signal
176    in contralateral motor cortex and operculum as well as ipsilateral cerebellum when contrasting hand
177    responses against each other (Fig. 6C, right panel). These results are in line with previous results on
178    outcome processing and response selection and thus assure the general data quality.

## fMRI: Biased learning in prefrontal cortex and striatum

179

180    To test which brain regions were involved in biased learning, we performed a model-based
181    GLM featuring the trial-by-trial PE update as a parametric regressor (see GLM notation in S06). We
182    used the group-level parameters of the best fitting computational model (M5) to compute trial-by-
183    trial belief updates (i.e., prediction error * learning rate) for every participant. In assessing neural
184    signatures of biased learning, we faced the complication that standard (Rescorla-Wagner learning in
185    M1) and biased PEs (winning model M5) are highly correlated. A mean correlation of 0.92 across
186    participants (range 0.88–0.95) made it difficult to neurally distinguish biased from standard learning.
187    To circumvent this collinearity problem, we decomposed the biased PE (computed using model M5)
188    into the standard PE (computed using model M1) plus a difference term (19, 34):

$$PE_{BIAS} = PE_{STD} + PE_{DIF}$$

189    A neural signature of biased learning should, significantly and with the same sign, encode
190    both components of this biased PE term. Standard PEs and difference term were uncorrelated (mean
191    correlation of -0.02 across participants; range -0.33–0.24). We tested for biased PEs $PE_{BIAS}$ by
192    computing which regions significantly encoded the conjunction of both its components, i.e., standard
193    prediction errors $PE_{STD}$ and the difference to biased PEs $PE_{DIF}$. While $PE_{STD}$ was encoded in a range of
194    cortical and subcortical regions (Fig. 3B, S07) previously reported in the literature (35), significant
195    encoding of both $PE_{STD}$ and PE$_{DIF}$ (conjunction) occurred in striatum (caudate, nucleus accumbens),
196    vmPFC/ perigenual ACC (area 32d), ventral ACC (area 23/24), posterior cingulate cortex (PCC), left
197    motor cortex, left inferior temporal gyrus, and early visual regions (Fig. 3C; see full report of
198    significant clusters in S07). Thus, BOLD signal in these regions was better described (i.e., more
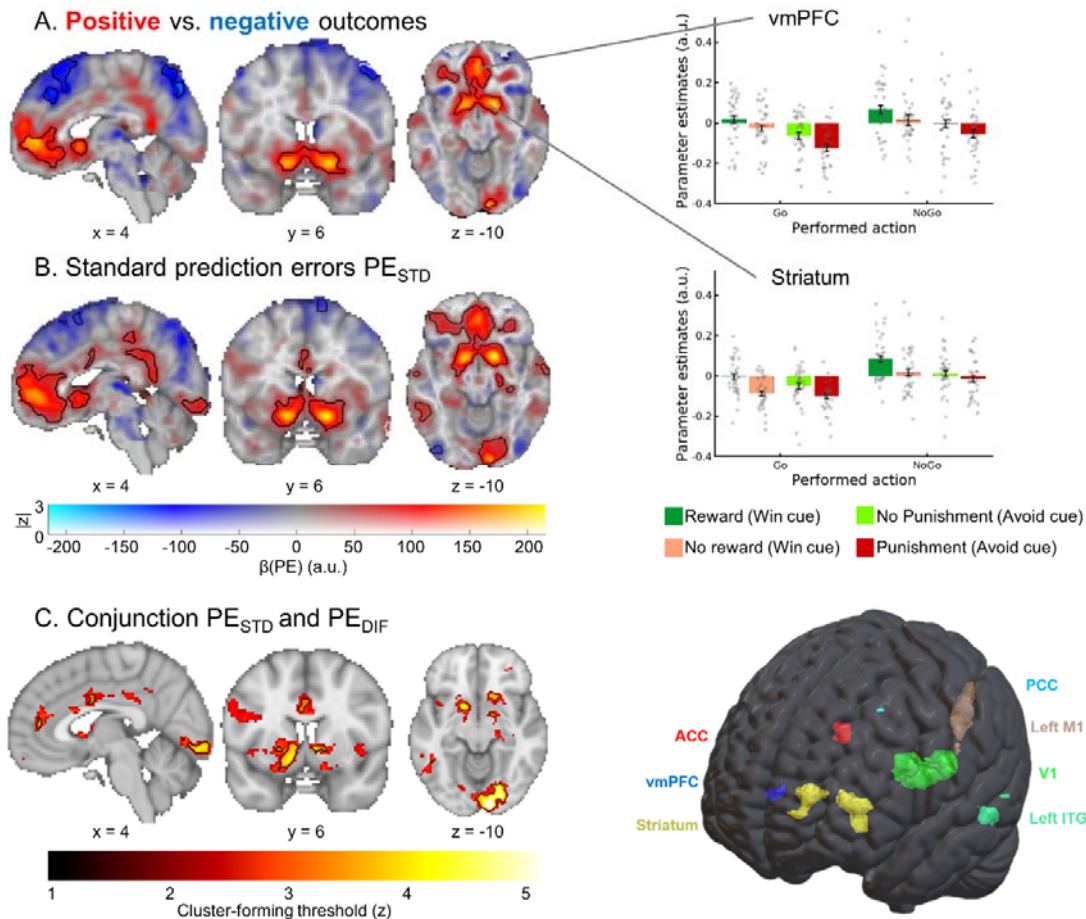199    variance explained) by biased learning than by standard prediction error learning.
200

**Figure 3. BOLD signal reflecting outcome processing**. BOLD effects displayed using a dual-coding visualization: color indicates the parameter estimates and opacity the associated z-statistics. Significant clusters are surrounded by black edges. (A) significantly higher BOLD signal for favorable outcomes (rewards, no punishments) compared with unfavorable outcomes (no rewards, punishments) was present in a range of regions including bilateral ventral striatum and vmPFC. Bar plots show mean parameter estimates per condition (±SEM across participants; dots indicating individual participants) (B) BOLD signals correlated positively to "standard" RL prediction errors in several regions, including the ventral striatum, vmPFC, PCC and ACC. (C) Left panel: Regions encoding both the standard PE term and the difference term to biased PEs (conjunction) at different cluster-forming thresholds (1 < z < 5, color coding; opacity constant). Clusters significant at a threshold of z > 3.1 are surrounded by black edges. In bilateral striatum, ACC, vmPFC, PCC, left motor cortex, left inferior temporal gyrus, and primary visual cortex, BOLD is significantly better explained by biased learning than by standard learning. Right panel: 3D representation with all seven regions encoding biased learning (used in fMRI-informed EEG analyses).

201

## EEG: Biased learning in midfrontal delta, theta, and beta power

203       Similar to the fMRI analyses, we next tested whether midfrontal power encoded biased PEs
204 rather than standard PEs. While fMRI provides spatial specificity of where PEs are encoded, EEG
205 power provides temporal specificity of when signals encoding prediction errors occur (*26, 31*). In line
206 with our fMRI analysis, we used the standard PE term $PE_{STD}$ and the difference to the biased PE
207 term $PE_{DIF}$ as trial-by-trial regressors for EEG power at each channel-time-frequency bin for each
208 participant and then performed cluster-based permutation tests across the *b*-maps of all
209 participants. Note that differently from BOLD signal, EEG signatures of learning typically do not

210   encode the full prediction error. Instead, PE sign (favorable vs. unfavorable outcomes) and PE
211   magnitude (saliency, surprise) have been found encoded separately in the theta and delta band,
212   respectively (*28–30*). We thus added PE sign as an additional regressor to test for separate correlates
213   of PE sign and PE magnitude. Note that PE sign is identical for standard and biased PEs; only PE
214   magnitude distinguishes both learning models.

215       Both midfrontal theta and beta power reflected PE sign: Theta power was higher for
216   unfavorable than favorable outcomes (225–475 ms, $p$ = .006; Fig. 4A-B), while beta power was higher
217   for favorable than unfavorable outcomes (300–1,250 ms, $p$ = .002; Fig. 4A, C). Differences in theta
218   power were clearly strongest over frontal channels, while the effect in the beta range was more
219   diffuse, spreading over frontal and parietal channels (Fig. 4B-C). All results held when the condition-
220   wise ERP was removed from the data (see S08), suggesting that differences between conditions were
221   due to induced (rather than evoked) activity (for results in the time domain, see S09).

222       Delta power was indeed positively, though not significantly correlated with both $PE_{STD}$ ($p$ =
223   0.074, Fig. 4E) and $PE_{DIF}$ ($p$ = 0.185; Fig. 4F). Only the sum of both terms, i.e., the $PE_{BIAS}$ term, was
224   significantly encoded by delta power (225–475 ms; $p$ = .017; Fig. 4D). For a similar observation in the
225   time-domain EEG signal, see S10. Beyond delta power, beta power correlated positively, though not
226   significantly with $PE_{STD}$ ($p$ = 0.110, Fig. 4E) and significantly negatively with $PE_{DIF}$ ($p$ = .001, 425 –
227   850 ms). Encoding of $PE_{BIAS}$ was not significant either ($p$ = 0.550, Fig 4D).

228       In sum, both midfrontal theta power (negatively) and beta power (positively) encoded PE
229   sign. In addition, delta power encoded PE magnitude (positively). This encoding was only significant
230   for biased PEs, but not standard PEs. Taken together, as was the case for BOLD signal, midfrontal EEG
231   power also reflected biased learning. As a next step, we tested whether the identified EEG
232   phenomena were correlated with trial-by-trial BOLD signal in identified regions. Crucially, this
233   allowed us to test whether EEG correlates of cortical learning precede EEG correlates of subcortical
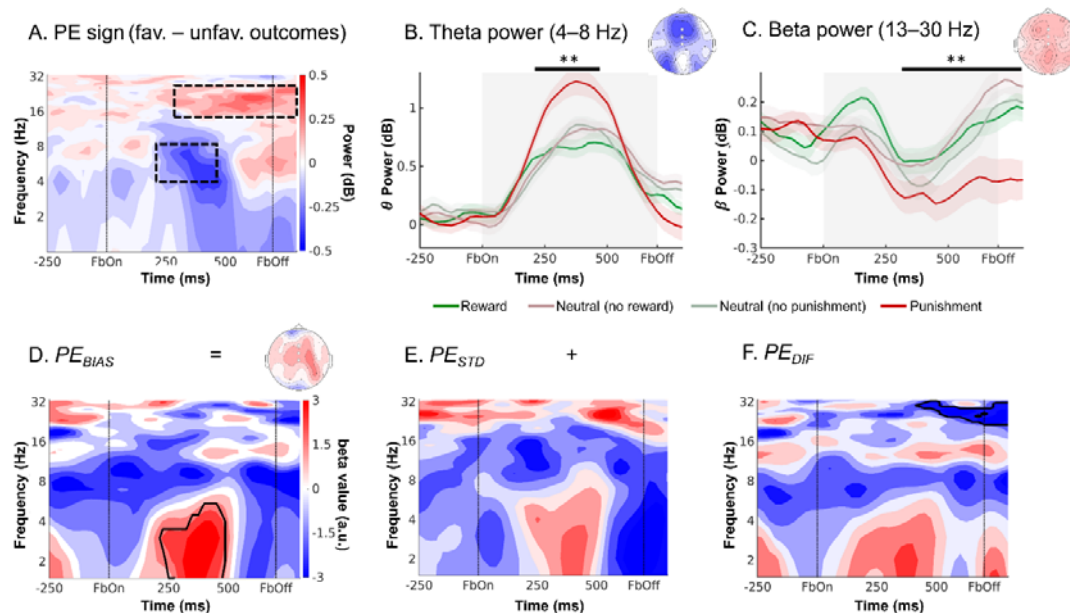234   learning.



**Figure 4. EEG time-frequency power over midfrontal electrodes (Fz/ FCz/ Cz). reflecting outcome processing.** (A) Time-frequency plot (logarithmic y-axis) displaying higher theta (4–8 Hz) power for unfavorable outcomes and higher beta power (16–32 Hz) for favorable outcomes. Black square dot boxes indicate clusters above threshold that drive significance in a-priori defined frequency ranges. (B). Theta power transiently increases for any outcome, but more so for unfavorable outcomes (especially punishments) around 225–475 ms after feedback onset. Black horizontal lines

indicate the time range for which the cluster driving significance is above threshold. (C) Beta power is higher for favorable than unfavorable outcomes over a long time period around 300–1,250 ms after feedback onset. (D-F). Correlations between midfrontal EEG power and trial-by-trial PEs controlling for PE sign. Solid black lines indicate clusters above threshold. Biased PEs were significantly positively correlated with midfrontal delta power (D). The correlations of delta with the standard PEs (E) and the difference term to biased PEs (F) were positive as well, though not significant. Beta power only encoded the difference term to biased PEs (F). ** $p$ < 0.01.

235

## Combined EEG-fMRI: Prefrontal cortex signals precede striatum during biased outcome processing

236
237

238      The observation that also cortical areas (vmPFC, ACC, PCC) show biased PEs is consistent with the
239  "external model" of cortical signals biasing learning processes in the striatum. However, this model
240  makes the crucial prediction that these bias signals should be present first in cortical areas and only
241  later in the striatum. Next, we used trial-by-trial BOLD signal from those regions encoding biased PE
242  to predict midfrontal EEG power. By determining the time points at which different regions
243  correlated with EEG power, we were able to infer the relative order of biased PE processing across
244  cortical and subcortical regions, revealing whether cortical processing preceded striatal processing.
245  We used trial-by-trial BOLD signal from the seven regions encoding biased PEs, i.e., striatum, ACC,
246  left motor cortex, vmPFC, PCC, left ITG, and primary visual cortex (see masks in S05) as regressors on
247  average EEG power over midfrontal electrodes (Fz/ FCz/ Cz). We controlled for biased PEs themselves
248  to capture additional variance in EEG explained by BOLD signal beyond the task regressors. As the
249  timeseries of all seven regions were included in one single regression, their regression weights reflect
250  each region's unique contribution, controlling for any shared variance. In line with the "external
251  model", BOLD signal from prefrontal cortical regions correlated with midfrontal EEG power earlier
252  after outcome onset than did striatal BOLD signal:
253      First, ACC BOLD was significantly negatively correlated with alpha/ theta power early after
254  outcome onset (100–575 ms, 2 – 17 Hz, $p$ = .016; Fig. 5A). This cluster started in the alpha/ theta
255  range and then spread into the theta/delta range (henceforth called "lower alpha band power"). It
256  was not observed in the EEG-only analyses reported above.
257      Second, while vmPFC/ perigenual ACC BOLD did not correlate significantly with midfrontal EEG
258  power (p = .184), BOLD in PCC was negatively correlated with theta/ delta power (Fig. 5B; 175–500
259  ms, 1–6 Hz, $p$ = .014). This finding bears resemblance in terms of time-frequency space to the cluster
260  of (negative) PE sign encoding in the theta band and (positive) PE magnitude encoding in the delta
261  band identified in the EEG-only analyses (Fig. 4A). As a reverse check of this link, we added the trial-
262  by-trial power in the EEG-only theta/delta band cluster as a regressor to the fMRI GLM featuring
263  prediction errors, which yielded significant clusters of negative EEG-BOLD correlation in vmPFC and
264  PCC (Fig. 5F; S13). We thus discuss vmPFC and PCC together in the following.
265      Third, there was a significant positive correlation between striatal BOLD and midfrontal beta/
266  alpha power (driven by a cluster at 100–800 ms, 7–23 Hz, $p$ = .010; Fig. 5C). This finding bears
267  resemblance in time-frequency space to the cluster of positive PE sign encoding in beta power
268  identified in the EEG-only analyses (Fig. 4A). Again, to substantiate this link, we performed the
269  reverse approach of using trial-by-trial power in the EEG-only beta band cluster as a regressor added
270  to the fMRI GLM. Clusters of positive EEG-BOLD correlations in right dorsal caudate (and left
271  parahippocampal gyrus) as well as clusters of negative correlations in bilateral dorsolateral PFC
272  (dlPFC) and supramarginal gyrus (SMG; Fig. 5G; see S13) confirmed the positive striatal BOLD-beta

273    power association. Given that the striatum is unlikely to be the source of midfrontal beta power over
274    the scalp, this analysis suggests dlPFC and SMG as likely candidate sources.
275        Finally, regarding the other three regions that showed a significant BOLD signature of biased PEs:
276    BOLD in left motor cortex was significantly negatively correlated with early midfrontal beta power ($p$
277    = .002; around 0 – 625 ms; see S11). There were no significant correlations between midfrontal EEG
278    power and left inferior temporal gyrus or primary visual cortex BOLD (see S11). All results were
279    robust to different analysis approaches including shorter trial windows, different GLM specifications,
280    inclusion of task-condition and fMRI motion realignment regressors, and individual modelling of each
281    region, and were not reducible to phenomena in the time domain (see S12).
282        In sum, there were negative correlations between ACC BOLD and midfrontal lower alpha band
283    power early after outcome onset, negative correlations between PCC BOLD and midfrontal theta/
284    delta power at intermediate time points, and positive correlations between striatal BOLD and
285    midfrontal beta power at late time points (Fig. 5D, H). These results are consistent with an "external
286    model" of motivational biases arising from early cortical processes biasing later learning processes in
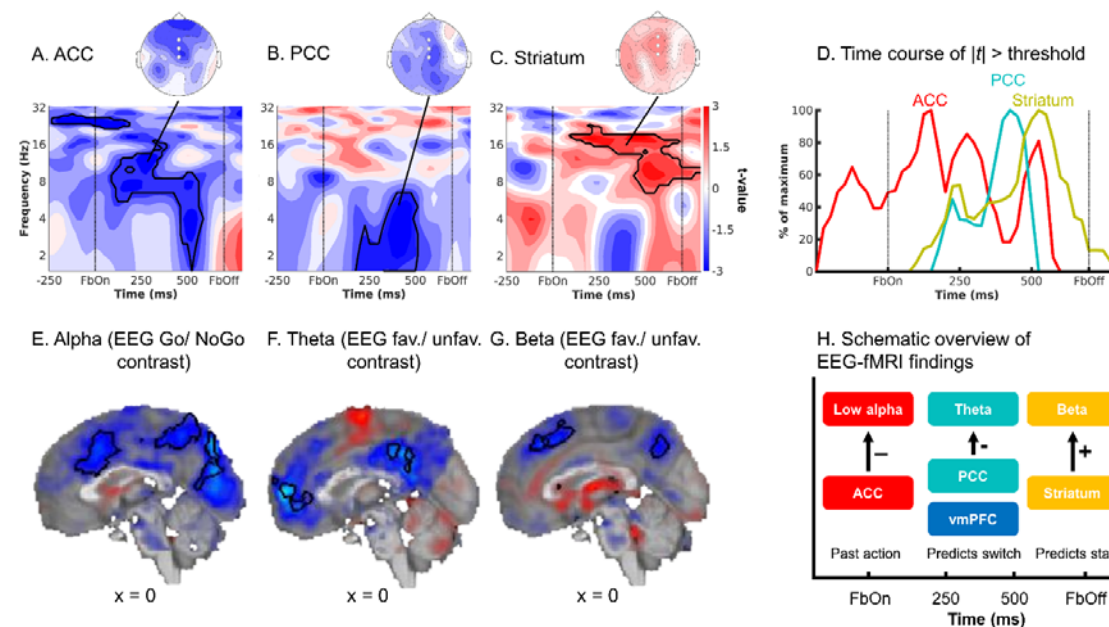287    the striatum.



**Figure 5. fMRI-informed EEG analyses**. Unique temporal contributions of BOLD signal in (A) ACC, (B) PCC, and (C) striatum to average EEG power over midfrontal electrodes (Fz/ FCz/ Cz). Group-level $t$-maps display the modulation of the EEG power by trial-by-trial BOLD signal in the selected ROIs. ACC BOLD correlates negatively with early alpha/ theta power, PCC BOLD negatively with theta/ delta power, striatal BOLD positively with beta/ alpha power. Areas surrounded by a black edge indicate clusters of $|t| > 2$ with $p < .05$ (cluster-corrected). Topoplots indicate the topography of the respective cluster. (D) Time course of ACC, PCC, and striatal BOLD correlations, normalized to the peak of the time course of each region. ACC-lower alpha band correlations emerge first, followed by (negative) PCC-theta correlations and finally positive striatum-beta correlations. Reverse approach using lower alpha (E), theta (F) and beta (G) power as trial-by-trial regressors in fMRI GLMs. These EEG-informed fMRI analyses corroborate the fMRI-informed EEG analyses: Lower alpha band power correlated negatively with the ACC BOLD, theta power negatively with vmPFC and PCC BOLD, and beta power positively with striatal BOLD. (H) Schematic overview of the main EEG-fMRI results: ACC encodes the previously performed response and correlates with early midfrontal lower alpha band power. vmPFC/ PCC (correlated with theta power) and striatum (correlated with beta power) both encode outcome valence, but have opposite effects on subsequent behavior. Note that activity in these regions temporally overlaps; boxes are ordered in temporal precedence of peak activity.

288

## ACC BOLD and midfrontal lower alpha band power encode the previously performed action during outcome presentation

289
290
291     While the clusters of EEG-fMRI correlation in the theta/ delta and beta range matched the
292     clusters identified in EEG-only analyses, the cluster of negative correlations between ACC BOLD and
293     early midfrontal lower alpha band power was novel and did not match our expectations. Given that
294     these correlations arose very soon after outcome onset, we hypothesized that ACC BOLD and
295     midfrontal lower alpha band power might reflect a process occurring even before outcome onset,
296     such as the maintenance ("eligibility trace") of the previously performed response to which credit
297     may later be assigned. We therefore assessed whether information of the previous response was
298     present in ACC BOLD and in the lower alpha band around the time of outcome onset.
299     First, we tested for BOLD correlates of the previous response at the time of *outcomes* (eight
300     outcome-locked regressors for every Go/ NoGo x reward/ no reward/ no punishment/ punishment
301     combination) while controlling for motor-related signals at the time of the *response* (response-locked
302     regressors for left-hand and right-hand button presses). At the time of outcomes, there was higher
303     BOLD signal for NoGo than Go responses across several cortical and subcortical regions, peaking in
304     both the ACC and striatum (Fig. 6E). This inversion of effects—higher BOLD for Go than NoGo
305     responses at the time of response (see quality checks), but the reverse at the time of outcome—was
306     also observed in the upsampled raw BOLD and was independent of the response of the next trial
307     (S14). In sum, large parts of cortex, including the ACC, indeed encoded the previously performed
308     response at the moment outcomes were presented, in line with the idea that the ACC maintains an
309     "eligibility trace" of the previously performed response.
310     Second, we tested for differences between Go and NoGo responses at the time of outcomes
311     in midfrontal broadband EEG power. Power was significantly higher on trials with Go than on trials
312     with NoGo responses, driven by clusters in the lower alpha band (spreading into the theta band;
313     around 0.000–0.425 sec., 1–11 Hz, $p$ = .012) and in the beta band (around 0.200–0.450 sec., 18–27
314     Hz, $p$ = .022; Fig. 6A, B). The first cluster matched the time-frequency pattern of ACC BOLD-alpha
315     power correlations (Fig. 5A).
316     If this activity cluster contained a signature of the previously performed response, it might have
317     been present throughout the delay between cue offset and outcome onset. When repeating the
318     above permutation test including the last second before outcome onset, there were significant
319     differences again, driven by a sustained cluster in the beta band (-1–0 sec., 13–33 Hz, $p$ = .002) and
320     two clusters in the alpha/ theta band (Cluster 1: -1.000– -0.275 sec., 1–10 Hz, $p$ = 0.014; Cluster 2: -
321     0.225–0.425 sec., 1–11 Hz, $p$ = .022; Fig. 6B). These findings suggest that lower alpha band power
322     might reflect a sustained memory of the previously performed response. Supplemental analyses
323     (S14) yielded that this Go-NoGo trace during outcome processing did not change over the time
324     course of the experiment, suggesting that it did not reflect typical fatigue/ time-on task effects often
325     observed in the alpha band.
326     Again, we performed the reverse EEG-fMRI analysis using trial-by-trial power in the identified
327     lower alpha band cluster (Fig. 6B) as an additional regressor in the quality-check fMRI GLM. Clusters
328     of negative EEG-BOLD occurred correlation in a range of cortical regions, including ACC and
329     precuneous (Fig. 5E; see S13). In sum, both ACC BOLD signal and midfrontal lower alpha band power
330     contained information about the previously performed response, consistent with the idea that both
331     signals reflect an "eligibility trace" of the response to which credit is assigned once an outcome is
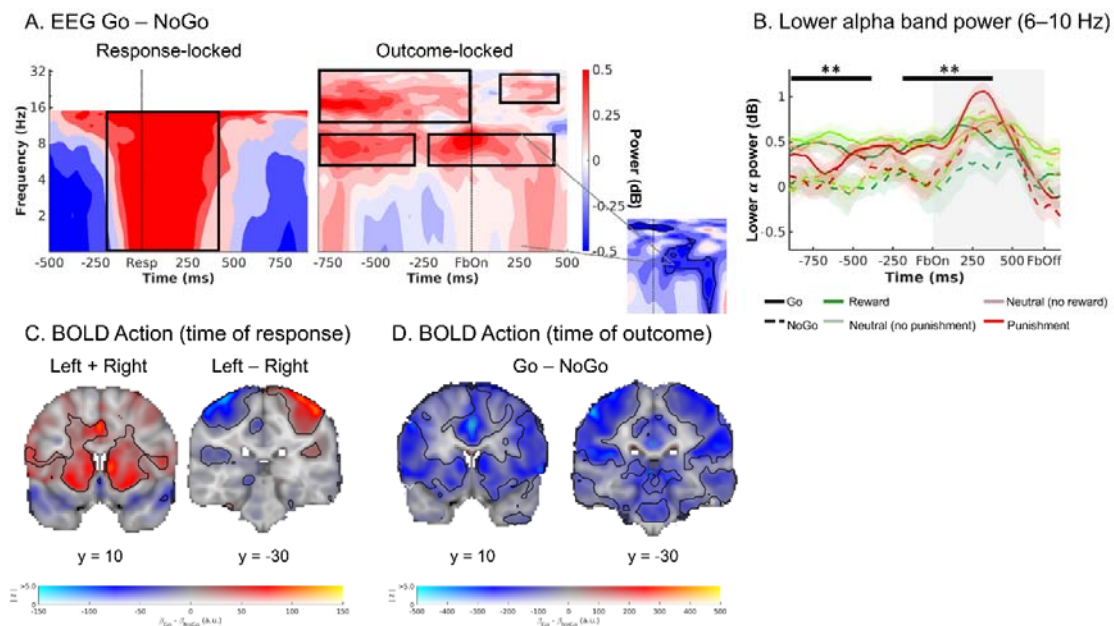332     obtained.

**Figure 6. Exploratory follow-up analyses on ACC BOLD signal and midfrontal lower alpha band power.** (A) Midfrontal time-frequency response-locked (left panel) and outcome-locked (right panel). Before and shortly after outcome onset, power in the lower alpha band is higher on trials with Go actions than on trials with NoGo actions. The shape of this difference resembles the shape of ACC BOLD-EEG TF correlations (small plot; note that this plot depicts BOLD-EEG correlations, which are negative). Note that differences between Go and NoGo trials occurred already before outcome onset in the alpha and beta range, reminiscent of delay activity, but were not fully sustained throughout the delay between response and outcome. (B) Midfrontal power in the lower alpha band per action x outcome condition. Lower alpha band power is consistently higher on trials with Go actions than on trials with NoGo actions, starting already before outcome onset. (C) BOLD signal differences between Go and NoGo actions (left panel) and left vs. right hand responses (right panel) at the time or responses. Response-locked ACC BOLD is significantly higher for Go than NoGo actions. (D) BOLD signal differences between Go and NoGo actions at the time of outcomes. Outcome-locked ACC BOLD (and BOLD in other parts of cortex) is significantly lower on trials with Go than on trials with NoGo actions.

## Striatal and vmPFC/ PCC BOLD differentially relate to action policy updating

EEG correlates of PCC BOLD and striatal BOLD occurred later than for the ACC BOLD, and overlapped with classical feedback-related midfrontal theta and beta power responses. We hypothesized that those neural signals might be more closely related to updating of action policies (i.e., which action to perform for each cue) and might thus predict the next response to the same cue (27, 36). We thus used the trial-by-trial BOLD responses in ACC, PCC, vmPFC and striatum to predict whether participants would repeat the same response on the next trial with the same cue ("stay") or switch to another response ("shift"). Mixed-effects logistic regression yielded that ACC BOLD did not significantly predict response repetition ($b$ = -0.019, $SE$ = 0.016, $\chi^2$(1) = 1.294, $p$ = .255). In contrast, BOLD in PCC/ vmPFC and striatum did predict response repetition, though in opposite directions: Participants were significantly more likely to repeat the same response when striatal BOLD was high ($b$ = 0.067, $SE$ = 0.024, $\chi^2$(1) = 9.051, $p$ = .003), but more likely to switch to another response when vmPFC BOLD ($b$ = -0.076, $SE$ = 0.017, $\chi^2$(1) = 15.559, $p$ < .001) or PCC BOLD ($b$ = -0.036, $SE$ = 0.016, $\chi^2$(1) = 3.691, $p$ = .030; Fig. 5H) was high (for plots, see S15). We also inspected the raw upsampled HRF shapes per region per condition, confirming that differential relationships were not driven by differences in HRF shapes across regions.

350    We also tested whether trial-by-trial midfrontal lower alpha band, theta, or beta power
351    (within the clusters identified in the EEG-only analyses) predicted action policy updating. Participants
352    were significantly more likely to repeat the same response when beta power was high ($b$ = 0.145, $SE$
353    = 0.041, $\chi^2(1)$ = 11.886, $p$ < .001), but more likely to switch when theta power was high ($b$ = -0.099, $SE$
354    = 0.047, $\chi^2(1)$ = 4.179, $p$ = .041). Notably, unlike its BOLD correlate in ACC, lower alpha band power
355    did predict response repetition, with more repetition when alpha power was high ($b$ = .0.179, SE =
356    0.052, $\chi^2(1)$ = 10.711, $p$ = .001; for plots, see S15).
357    In sum, high striatal BOLD and midfrontal beta power predicted that the same response
358    would be repeated on the next encounter of a cue, while high vmPFC and PCC BOLD and high theta
359    power predicted that participants would switch to another response. Thus, although both striatal and
360    vmPFC/PCC BOLD positively encoded biased prediction errors, these two sets of regions had opposite
361    roles in learning: while the striatum reinforces previous responses, vmPFC/PCC trigger the shift to
362    another response strategy (Fig. 5H).
363

## Discussion

365    We investigated neural correlates of biased learning for Go and NoGo responses. In line with
366    previous research (3, 9), participants' behavior was best described by a computational model
367    featuring faster learning from rewarded Go responses and slower learning from punished NoGo
368    responses. Neural correlates of biased PEs were present in BOLD signals in several regions, including
369    ACC, PCC, vmPFC, and striatum. These regions exhibited distinct midfrontal EEG power correlates.
370    Most importantly, correlates of prefrontal cortical BOLD preceded correlates of striatal BOLD: Trial-
371    by-trial ACC BOLD correlated with lower alpha band power immediately after outcome onset,
372    followed by PCC (and vmPFC) BOLD correlated with theta power, and finally striatal BOLD correlated
373    with beta power. These results are in line with a model of PFC biasing striatal outcome processing,
374    giving rise to motivational learning biases in behavior.
375

**Biased learning in PFC precedes the striatum**

377    The dominant idea about the origin of motivational biases has been that these biases are an
378    emergent feature of the asymmetric direct/ indirect pathway architecture in the basal ganglia (2, 16).
379    We find evidence that these biases are present first in prefrontal cortical areas, notably ACC and
380    vmPFC. This argues against biases purely being a "fixed" leftover of evolutionary ancient, subcortical
381    circuits. Rather, motivational learning biases might be an instance of sophisticated, even "model-
382    based" learning processes in the striatum instructed by the prefrontal cortex (37, 38). An influence of
383    PFC on striatal RL has prominently been observed in the case of model-based vs. model-free learning
384    (20, 21) and has been stipulated as a mechanism of how instructions can impact RL learning (17, 18).
385    Although there are reports of striatal processes preceding prefrontal processes within learning tasks
386    (39, 40), the opposite pattern of PFC preceding striatum has been observed as well (41) and a causal
387    impact of PFC on striatal learning is well established (42, 43).
388    The particular subregion of PFC showing the earliest EEG correlates was the ACC. This
389    observation is in line with an earlier EEG-fMRI study reporting ACC to be part of an early valuation
390    system preceding a later system comprising vmPFC and striatum (44). The ACC has been suggested to
391    encode models of agents' environment (45, 46) that are relevant for interpreting outcomes. ACC
392    BOLD has been found to scale with the size of PEs (22, 23), indexing how much should be learned
393    from new outcomes. We hypothesize that, at the moment of outcome, ACC maintains an "eligibility
394    trace" of the previously performed response (47), which might modulate the processing of outcomes

395 as soon as they become available (48, 49). Notably, ACC exhibited stronger BOLD signal for Go than
396 NoGo responses at the time of participants' response, but this pattern reversed at the time of
397 outcomes. This reversal rules out the possibility that response-locked BOLD signal simply spilled over
398 into the time of outcomes. Future research will be necessary to corroborate such a motor "eligibility
399 trace" in ACC.
400 In sum, the ACC might be in a designated position to inform subsequent outcome processing
401 in downstream regions by modulating the learning rate as a function of previously performed
402 response and the obtained outcome. Rather than striatal circuits being sufficient for the emergence
403 of motivational biases, the more "flexible" PFC seems to play role in instructing downstream striatal
404 learning processes.
405
406 **Striatum and midfrontal beta power signal maintenance of action policies**
407 Striatal, vmPFC and PCC BOLD encoded biased PEs. In line with previous research, striatal
408 BOLD positively linked to midfrontal beta power (50, 51), which positively encoded PE sign (25, 31,
409 52). PCC and vmPFC BOLD negatively linked to midfrontal theta/ delta power (32, 53, 54), which
410 encoded PE sign negatively, but PE magnitude positively. Notably, theta/ delta power correlates of
411 vmPFC/ PCC BOLD preceded beta power correlates of striatal BOLD in time, which aligns with
412 previous findings of motivational response biases being first visible in the vmPFC BOLD before they
413 impact striatal action selection (32).
414 Positive encoding of prediction errors in striatal BOLD signal is a well-established phenomenon
415 (35, 55). Striatal BOLD was better described by biased PEs than by standard PEs, corroborating the
416 presence of motivational learning biases also in striatal learning processes. Notably, EEG correlates of
417 striatal BOLD peaked rather late, suggesting that these processes are informed by early sources in
418 PFC which are connected to the striatum via recurrent feedback loops (15, 56). Positive prediction
419 errors increase the value of a performed action and thus strengthen action policies. Hence, it is not
420 surprising that high striatal BOLD signal and midfrontal beta power predicted action repetition (57,
421 58).
422
423 **vmPFC and midfrontal theta/ delta power signal updating of action policies**
424 In contrast to striatal learning signals, the PCC and vmPFC BOLD as well as midfrontal theta
425 and delta power signals were more complicated: Theta encoded PE sign, delta encoded PE
426 magnitude. Both correlates showed opposite polarities. This observation is in line with previous
427 literature suggesting that midfrontal theta and delta power (resp. the feedback-related negativity
428 and reward positivity components in the time domain EEG signal) might reflect the "saliency" or
429 "surprise" aspect of PEs (28, 29, 59). Surprises have the potential to disrupt an ongoing action policy
430 (60) and motivate a shift to another policy, which might explain why these signals predicted
431 switching to another response (61, 62). Notably, this EEG surprise signal was only significantly
432 correlated with the biased (but not the standard) PE term, corroborating that the surprise attributed
433 to outcomes depends on previously performed response, reflecting motivational learning biases. In
434 sum, both vmPFC and striatum encode biased PEs, though with different consequences for future
435 action policies.
436
437 **Limitations**
438 Taken together, distinct brain regions processed outcomes in a biased fashion at distinct time
439 points with distinct EEG power correlates. Simultaneous EEG-fMRI recordings allowed us to infer
440 when those regions reached their peak activity (63). However, the correlational nature of BOLD-EEG

441 links precludes strong statements about these regions actually generating the respective power
442 phenomena. Alternatively, activity in those regions might merely modulate the amplitude of time-
443 frequency responses originating from other sources. Furthermore, while the observed associations
444 align with previous literature (32, 50, 51, 53, 54), the considerable distance of the striatum to the
445 scalp raises the question whether scalp EEG could in principle reflect striatal activity, at all (64, 65).
446 Intracranial recordings have observed beta oscillations during outcome processing in the striatum
447 before (58, 66, 67). Also, our analysis controlled for BOLD signal in motor cortex, an alternative
448 candidate source for beta power, suggesting that late midfrontal beta power does not merely reflect
449 motor cortex beta. Even if the striatum is not the generator of the beta oscillations over the scalp,
450 their true (cortical) generator might be tightly coupled to the striatum and thus act as a "transmitter"
451 of striatal beta oscillations. In fact, the analyses using trial-by-trial beta power to predict BOLD
452 yielded significant clusters in dlPFC and SMG, two candidate regions for such a "transmitter".
453 Finally, the correlational nature of the study prevents strong statements over any causal
454 interactions between the observed regions. We assume here that a region showing an earlier
455 midfrontal EEG correlate influences other regions showing later midfrontal EEG correlates, and such
456 an influence is plausible given findings of feedback loops between prefrontal regions and the
457 striatum (56). Future studies targeting those regions via selective causal manipulations will be
458 necessary to test for the causal role of PFC in informing striatal learning.
459
460 **The role of motivational biases in credit assignment and learning**
461 In conclusion, biased learning—increased credit assignment to rewarded action, decreased credit
462 assignment to punished inaction—was visible both in behavior and in BOLD signal in a range of
463 regions. EEG correlates of prefrontal cortical regions, notably ACC and vmPFC, *preceded* correlates of
464 the striatum, consistent with a model of the PFC biasing RL in the striatum. The ACC appeared to hold
465 a "motor eligibility trace" of the past response, biasing early outcome processing. Subsequently,
466 biased learning was also present in vmPFC/ PCC and striatum, with opposite roles in adjusting vs.
467 maintaining action policies. These results refine previous views on the neural origin of these learning
468 biases, which might not purely be "naïve" remnants of evolutionary ancient, "primitive" parts of the
469 brain, but rather incorporate sophisticated, even "model-based" processes relying on frontal inputs.
470 The PFC is typically believed to facilitate goal-directed over instinctive processes. Hence, PFC
471 involvement into biased learning suggests that these biases are not necessarily agents' inescapable
472 "fate", but rather likely act as global "priors" that facilitate learning of more local relationships. They
473 allow for combining "the best of both worlds"—long-term experience with consequences of actions
474 and inactions together with flexible learning from rewards and punishments.

475 # Materials and methods

476 ## Participants
477 Thirty-six participants ($M_{age}$ = 23.6, $SD_{age}$ = 3.4, range 19–32; 25 women; all right-handed; all normal
478 or corrected-to-normal vision) took part in a single 3-h data collection session, for which they
479 received €30 flat fee plus a performance-dependent bonus (range €0–5, $M_{bonus}$ = €1.28, $SD_{bonus}$ =
480 1.54). The study was approved by the local ethics committee (CMO2014/288; Commissie
481 Mensengeboden Onderzoek Arnhem-Nijmegen) and all participants provided written informed
482 consent. Exclusion criteria comprised claustrophobia, allergy to gels used for EEG electrode
483 application, hearing aids, impaired vision, colorblindness, history of neurological or psychiatric
484 diseases (including heavy concussions and brain surgery), epilepsy and metal parts in the body, or
485 heart problems. Sample size was based on previous EEG studies with a comparable paradigm (9, 68).

486         Behavioral and modeling results include all 36 participants. The following participants were
487    excluded from analyses of neural data: For two participants, fMRI functional-to-standard image
488    registration failed; hence, all fMRI-only results are based on 34 participants ($M_{age}$ = 23.47, 25
489    women). Four participants exhibited excessive residual noise in their EEG data (> 33% rejected trials)
490    and were thus excluded from all EEG analyses; hence, all EEG-only analyses are based on 32
491    participants ($M_{age}$ = 23.09, 23 women). For combined EEG-fMRI analyses, we excluded the above-
492    mentioned six participants plus one more participant whose regression weights for every regressor
493    were about ten times larger than for other participants, leaving 29 participants ($M_{age}$ = 23.00, 22
494    women). Exclusions were in line with a previous analysis of this data set (32). fMRI- and EEG-only
495    results held when analyzing only those 29 participants (see S01).

## Task

497         Participants performed a motivational Go/ NoGo learning task (3, 9) administered via
498    MATLAB R2014b (MathWorks, Natick, MA, United States) and Psychtoolbox-3.0.13. On each trial,
499    participants saw a gem-shaped cue for 1300 ms, which signaled whether they could potentially win a
500    reward (Win cues) or avoid a punishment (Avoid cues), and whether they had to perform a Go (Go
501    cue) or NoGo response (NoGo cue). They could press a left ($Go_{LEFT}$), right ($Go_{RIGHT}$), or no (NoGo)
502    button while the cue was presented. Only one response option was correct per cue. Participants had
503    to learn both cue valence and required action from trial-and-error. After a variable inter-stimulus-
504    interval of 1,400–1,600 ms, the outcome was presented for 750 ms. Potential outcomes were a
505    reward (symbolized by coins falling into a can) or neutral outcome (can without money) for Win cues,
506    and a neutral outcome or punishment (symbolized by money falling out of a can) for Avoid cues.
507    Feedback validity was 80%, i.e., correct responses were followed by favorable outcomes (rewards/
508    no punishments) on only 80% of trials, while incorrect responses were still followed by favorable
509    outcomes on 20% of trials. Trials ended with a jittered inter-trial interval of 1250–2000 ms, yielding
510    total trial lengths of 4700–6650 ms.
511         Participants gave left and right Go responses via two button boxes positioned lateral to their
512    body. Each box featured four buttons, but only one button per box was required in this task. When
513    participants accidentally pressed a non-instructed button, they received the message "Please press
514    one of the correct keys" instead of an outcome. In the analyses, these responses were recoded into
515    the instructed button on the respective button box. In the fMRI GLMs, such trials were modeled with
516    a separate regressor.
517         Before the task, participants were instructed that each cue could be followed by either
518    reward or punishment, each cue had one optimal response, feedback was probabilistic, and that the
519    rewards and punishments were converted into a monetary bonus upon completion of the study.
520    They performed an elaborate practice session in which they got familiarized first with each condition
521    separately (using practice stimuli) and finally practiced all conditions together. They then performed
522    640 trials of the main task, separated into two sessions of 320 trials with separate cue sets.
523    Introducing a new set of cues allowed us to prevent ceiling effects in performance and investigate
524    continuous learning throughout the task. Each session featured eight cues that were presented 40
525    times. After every 100–110 trials (~ 6 min.), participants could take a self-paced break. The
526    assignment of the gems to cue conditions was counterbalanced across participants, and trial order
527    was pseudo-random (preventing that the same cue occurred on more than two consecutive trials).

## Behavior analyses

We used mixed-effects logistic regression (as implemented in the R package *lme4*) to analyze behavioral responses (Go vs. NoGo) as a function of required action (Go/ NoGo), cue valence (Win/ Avoid), and their interaction. We included a random intercept and all possible random slopes and correlations per participant to achieve a maximal random-effects structure (*69*). Sum-to-zero coding was employed for the factors. Type 3 *p*-values were based on likelihood ratio tests (implemented in the R package *afex*). We used a significance criterion of α = .05 for all the analyses.

Furthermore, we used mixed-effects logistic regression to analyze "stay behavior", i.e., whether participants repeated an action on the next encounter of the same cue, as a function of outcome valence (positive: reward or no punishment/ negative: no reward or punishment), outcome salience (salient: reward or punishment/ neutral: no reward or no punishment), and performed action (Go/ NoGo). We again included all possible random intercepts, slopes, and interactions.

## Computational modeling

We fit a series of increasingly complex RL models to participants' choices to decide between different algorithmic explanations for the emergence of motivational biases in behavior. We employed the same set of nested models as in previous studies using this task (*3, 9*). For tests of alternative biases specifications, see S03.

### Model space

To determine whether a Pavlovian response bias, an instrumental learning bias, or both biases jointly predicted behavior best, we fitted a series of increasing complex computational models. In each trial (t), choice probabilities for all three response options (a) given the displayed cue (s) were computed from their action weights (modified Q-values) using a softmax function:

$$p(a_t|s_t) = \frac{\exp\,(w(a_t, s_t))}{\sum_a \exp\,(w(a\prime, s_t))} \tag{1}$$

After each response, action values were updated with the prediction error based on the obtained outcome $r \in \{-1; 0; 1\}$. As the starting model (M1), we fitted an standard delta-learning model (*70*) in which action values were updated with prediction errors, i.e.,the deviation between the experienced outcome and expected outcome. This model contained two free parameters: the learning rate (ε) scaling the updating term and the feedback sensitivity (ρ) scaling the received outcome:

$$Q_t(a_t, s_t) = Q_{t-1}(a_t, s_t) + \varepsilon(\rho r - Q_{t-1}(a_t, s_t)) \tag{2}$$

In this model, choice probabilities were fully determined by action values, without any bias. We assigned cue valence $V_s$ to 0.5 for Win cues and -0.5 for Avoid cues and used cue valence scaled by participants' individual feedback sensitivity as initial action values $Q_0$. Unlike previous versions of the task (*3*), cue valences were not instructed, but had to be learned from outcomes, as well (*9*). Thus, until experiencing the first reward/ punishment for a cue, participants could not know its valence (and not learn from neutral feedback). Hence, for these trials, action values were multiplied with zero when computing choice probabilities. After the first encounter of a valenced outcome, action values were "unmuted" and started to influence choices probabilities, retrospectively considering all previous outcomes.

In M2, we added the Go bias parameter $b$, which accounted for individual differences in participants' overall propensity to make Go responses, to the action values Q, resulting in action weights w:

$$w(a_t, s_t) = \begin{cases} Q_t(a_t, s_t) + b & if\ a = Go \\ Q_t(a_t, s_t) & else \end{cases} \tag{3}$$

571        In M3, we added a Pavlovian response bias $\pi$, scaling how positive/ negative cue valence

572    (Pavlovian values) increased/ decreased the weights of Go responses:

573
$$w(a_t, s_t) = \begin{cases} Q_t(a_t, s_t) + b + \pi V(s) & if\ a = Go \\ Q_t(a_t, s_t) & else \end{cases} \quad (4)$$

574        We assigned cue valence to 0.5 for Win cues and -0.5 for Avoid cues. Cue valence became

575    effective only once the participant had experienced the first reward/ punishment for that cue;

576    beforehand, it was treated as zero. The Pavlovian response bias affected left-hand and right-hand Go

577    responses similarly and thus reflected generalized activation/ inactivation by the cue valence.

578        In M4, we added an instrumental learning bias $\kappa$, increasing the learning rate for rewards

579    after Go responses and decreasing it for punishments after NoGo responses:

580
$$\varepsilon = \begin{cases} \varepsilon_0 + \kappa & if\ r_t = 1\ and\ a = go \\ \varepsilon_0 - \kappa & if\ r_t = -1\ and\ a = nogo \\ \varepsilon_0 & else \end{cases} \quad (5)$$

581        The instrumental learning bias was specific to the response shown, thus reflecting a specific

582    enhancement in action learning/ impairment in unlearning for that particular response.

583        In the model M5, we included both the Pavlovian response bias and the instrumental

584    learning bias.

585        The hyperpriors were $X_\rho \sim \mathcal{N}(2,3)$, $X_\varepsilon \sim \mathcal{N}(0,2)$, $X_{b,\pi,\kappa} \sim \mathcal{N}(0,3)$. For computing the

586    participant-level parameters, $\rho$ was exponentiated to constrain it to positive values, and the inverse-

587    logit transformation was applied to $\varepsilon$ to constraint it to the range [0 1]. We made sure that the effect

588    of $\kappa$ on $\varepsilon$ was symmetrical by computing it as:

589
$$\varepsilon = \begin{cases} \varepsilon_0 = inv.\,logit(\varepsilon) \\ \varepsilon_{punished\ NoGo} = inv.\,logit(\varepsilon - \kappa) & if\ \varepsilon_0 < .5 \\ \varepsilon_{rewarded\ Go} = \varepsilon_0 + (\varepsilon_0 - \varepsilon_{punished\ NoGo}) & if\ \varepsilon_0 < .5 \end{cases} \quad (6)$$
$$\varepsilon = \begin{cases} \varepsilon_{rewarded\ Go} = inv.\,logit(\varepsilon + \kappa) & if\ \varepsilon_0 > .5 \\ \varepsilon_{punished\ NoGo} = \varepsilon_0 + (\varepsilon_0 - \varepsilon_{rewarded\ Go}) & if\ \varepsilon_0 > .5 \end{cases}$$

590    **Model fitting and comparison**

591        For model fitting and comparison, we used Hierarchical Bayesian inference as implemented

592    in the CBM toolbox in Matlab (71). This approach combines hierarchical Bayesian parameter

593    estimation with random-effects model comparison (72). The fitting procedure involves two steps,

594    starting with the Laplace approximation of the model evidence to compute the group evidence,

595    which quantifies how well each model fits the data while penalizing for model complexity. Both

596    group-level and individual-level parameters are estimated using an iterative algorithm. We used wide

597    Gaussian priors (see hyperpriors above) and exponential and sigmoid transforms to constrain

598    parameter spaces. Subsequent random-effects model selection allows for the possibility that

599    different models generated the data for different participants. Participants contribute to the group-

600    level parameter estimation in proportion to how well a given model fits their data, quantified via a

601    responsibility measure (i.e., the probability that the model at hand is responsible for generating data

602    of the respective participant). This model-comparison approach has been shown to be less

603    susceptible to the influence of outliers (71). We selected the "winning" model based on the

604    protected exceedance probability.

605    **Model validation**

606        We assured that the winning model was able to reproduce the data, using the sampled

607    combinations of participant-level parameter estimates to create 3,600 agents that "played" the task.

608     We employed two approaches to simulate the task: *posterior predictive model simulations* and *one-*
609     *step-ahead model predictions*. In the posterior predictive model simulations, agents' choices were
610     sampled probabilistically based on their action values, and outcomes probabilistically sampled based
611     on their choices. This method ignores participant-specific choice histories and can thus yield choice/
612     outcome sequences that diverge considerably from participants' actual experiences. In contrast, one-
613     step-ahead predictions use participants' actual choices and experienced outcomes in each trial to
614     update action values. We simulated choices for each participant using both methods, which
615     confirmed that the winning model M5 ("asymmetric pathways model") was able to qualitatively
616     reproduce the data, while an alternative implementation of biased learning ("action priming model")
617     failed to do so (see S03).

## fMRI data acquisition

619     fMRI data were collected on a 3T Siemens Magnetom Prisma fit MRI scanner with a 64-
620     channel head coil. During scanning, participants' heads were restricted using foam pillows and strips
621     of adhesive tape were applied to participants' forehead to provide active motion feedback and
622     minimize head movement (*73*). After two localizer scans to position slices, we collected functional
623     scans with a whole-brain T2*-weighted sequence (68 axial-oblique slices, TR = 1400 ms, TE = 32 ms,
624     voxel size 2.0 mm isotropic, interslice gap 0 mm, interleaved multiband slice acquisition with
625     acceleration factor 4, FOV 210 mm, flip angle 75°, A/ P phase encoding direction). The first seven
626     volumes of each run were automatically discarded. This sequence was chosen because of its balance
627     between a short TR and relatively high spatial resolution, which was required to disentangle cue and
628     outcome-related neural activity. Pilots using different sequences yielded that this sequence
629     performed best in reducing signal loss in striatum.
630     Furthermore, after task completion, we removed the EEG cap and collected a high-resolution
631     anatomical image using a T1-weighted MP-RAGE sequence (192 sagittal slices per slab, GRAPPA
632     acceleration factor = 2, TI = 1100 ms, TR = 2300 ms, TE = 3.03 ms, FOV 256 mm, voxel size 1.0 mm
633     isotropic, flip angle 8°) which was used to aid image registration, and a gradient fieldmap (GRE; TR =
634     614 ms, TE1 = 4.92 ms, voxel size 2.4 mm isotropic, flip angle 60°) for distortion correction. For one
635     participant, no fieldmap was collected due to time constraints. At the end of each session, an
636     additional DTI data collection took place; results will be reported elsewhere.

## fMRI preprocessing

638     All fMRI pre-processing was performed in FSL 6.0.0. After cleaning images from non-brain
639     tissue (brain-extraction with BET), we performed motion correction (MC-FLIRT), spatial smoothing
640     (FWHM 3 mm), and used fieldmaps for B0 unwarping and distortion correction in orbitofrontal areas.
641     We used ICA-AROMA (*74*) to automatically detect and reject independent components associated
642     with head motion. Finally, images were high-pass filtered at 100 s and pre-whitened. After the first-
643     level GLM analyses, we computed and applied co-registration of EPI images to high-resolution images
644     (linearly with FLIRT using boundary-based registration) and to MNI152 2mm isotropic standard space
645     (non-linearly with FNIRT using 12 DOF and 10 mm warp resolution).

## ROI selection

647     For fMRI-informed EEG analyses, we first created a functional mask as the conjunction of the
648     $PE_{STD}$ and $PE_{DIF}$ contrasts by thresholding both z-maps at z > 3.1, binarizing, and multiplying them (see
649     S05). After visual inspection of the respective clusters, we created seven anatomical masks based on
650     the probabilistic Harvard-Oxford Atlas (thresholded at 10%): striatum and ACC (see above), vmPFC

651  (combined frontal pole, frontal medial cortex, and paracingulate gyrus), motor cortex (combined
652  precentral and postcentral gyrus), PCC (Cingulate Gyrus, posterior division), ITG (Inferior Temporal
653  Gyrus, posterior division, and Inferior Temporal Gyrus, temporooccipital part) and primary visual
654  cortex (Lingual Gyrus, Occipital Fusiform Gyrus, Occipital Pole). We then multiplied this functional
655  mask with each of the seven anatomical masks, returning seven masks focused on the respective
656  significant clusters, which were then used for signal extraction. For the ACC mask, we manually
657  excluded voxels in subgenual ACC belonging to a distinct cluster. Masks were back-transformed to
658  each participant's native space.
659          For bar plots in Fig. 3A, we multiplied the anatomical masks of vmPFC and striatum specified
660  above with the binarized outcome valence contrast.

661  ## fMRI analyses
662          For each participant, data were modelled using two event-related GLMs. First, we performed
663  a model-based GLM in which used trial-by-trial estimates of biased PEs as regressors. Second, we
664  used another model-free GLM in which we model all possible action x outcome combinations via
665  outcome-locked categorical regressors while at the same time modeling response-locked left- and
666  right-hand response regressors. This model free GLM also contained the valence contrast reported as
667  an initial manipulation check.
668          In the model-based GLM, we used two model-based regressors that reflected the trial-by-
669  trial prediction error (PE) update term. For this purpose, we extracted the group-level parameters of
670  the best fitting computational model M5 (asymmetric pathways model) and used those parameters
671  to compute the prediction error on every trial for every participant. Using the same parameter for
672  each participant is warranted when testing for the same qualitative learning pattern across
673  participants (75). Given that both standard (base model M1) and biased (winning model M5) PEs
674  were highly correlated (mean correlation of 0.921 across participants, range 0.884–0.952), it
675  appeared difficult to distinguish standard learning from biased learning. As a remedy, we
676  decomposed the biased PE into the standard PE plus a difference term as $PE_{BIAS} = PE_{STD} + PE_{DIF}$
677  (19, 34). Any region displaying truly biased learning should significantly encode *both* the standard PE
678  term and the difference term. The standard PE and difference term were much less correlated (mean
679  correlation of -0.020, range -0.326–0.237). To control for cue-related activation, we furthermore
680  added four regressors spanned by crossing cue valence and performed action (Go response to Win
681  cue, Go response to Avoid cue, NoGo response to Win cue, NoGo response to Avoid cue).
682          The model-free GLM included a separate regressor for each of the eight conditions obtained
683  when crossing performed action (Go/ NoGo) and obtained outcome (reward/ no reward/ no
684  punishment/ punishment). We fitted four contrasts: 1) one contrast comparing conditions with
685  favorable (reward/ no punishment) and non-favorable (no reward/ punishment) outcomes, used as a
686  quality check to identify regions that encoded outcome valence; 2) one contrast comparing Go vs.
687  NoGo responses at the time of the outcome; 3) one contrast summing of left- and right-hand
688  responses, reflecting Go vs. NoGo responses at the time of the response; and 4) one contrast
689  subtracting right- from left-handed responses, reflecting lateralized motor activation. As this GLM
690  resulted in empty regressors for several participants when fitted on a block level, making it
691  impossible to use the data of the respective blocks on a higher level, we instead concatenated blocks
692  and performed a single GLM per participant. We therefore registered the data from all blocks to the
693  middle image of the first block (default reference volume in FSL) using MCFLIRT. The first and last 20

694     seconds of each block did not feature any task-related events, such that carry-over effects of task
695     events in the design matrix from one block to another were not possible.
696         In both GLMs, we added four regressors of no interest: one for the motor response (left = +1,
697     right = -1, NoGo = 0), one for error trials, one for outcome onset, and one for trials with invalid motor
698     response (and no outcome respectively). We also added nine or more nuisance regressors: the six
699     realignment parameters from motion correction, mean cerebrospinal fluid (CSF) signal, mean out-of-
700     brain (OBO) signal, and a separate spike regressor for each volume with a relative displacement of
701     more than 2 mm (occurred in 10 participants; in those participants: M = 7.40, range 1–29). For the
702     model-free GLM, nuisance regressors were added separately for each block as well as an overall
703     intercept per block. We convolved task regressors with double-gamma haemodynamic response
704     function (HRF) and high-pass filtered the design matrix at 100 s.
705         First-level contrasts were fit in native space. Afterwards, co-registration and reslicing was
706     applied to participants' contrast maps, which were then combined on a (participant and) group level
707     using FSL's mixed effects models tool FLAME with a cluster-forming threshold of z > 3.1 and cluster-
708     level error control at α < .05 (i.e., two one-sided tests with α < .025).

## EEG data acquisition

709
710         We recorded EEG data with 64 channels (BrainCap-MR-3-0 64Ch-Standard; Easycap GmbH;
711     Herrsching, Germany; international 10-20 layout, reference electrode at FCz) plus channels for
712     electrocardiogram, heart rate, and respiration (used for MR artifact correction) at a sampling rate of
713     1000 Hz. We placed MRI-compatible EEG amplifiers (BrainAmp MR plus; Brain Products GmbH,
714     Gilching, Germany) behind the MR scanner and attached cables to the participants once they were
715     located in final position in the scanner. Furthermore, we fixated cables using sand-filled pillows to
716     reduce artifacts induced through cable movement in the magnetic field. During functional scans, the
717     MR helium pump was switched off to reduce EEG artifacts. After the scanning, we recorded the exact
718     EEG electrode locations on participants' heads relative to three fiducial points using a Polhemus
719     FASTRAK device. For four participants, no such data were available due to time constraints/ technical
720     errors, in which case we used the average electrode locations of the remaining 32 participants.

## EEG pre-processing

721
722         First, raw EEG data were cleaned from MR scanner and cardioballistic artifacts using
723     BrainVisionAnalyzer (76). The rest of the pre-processing was performed in Fieldtrip (77). After
724     rejecting channels with high residual MR noise (mean 4.8 channels per participant, range 1–13), we
725     epoched trials into time windows of -1,400–2,000 ms relative to the onset of outcomes. Timing of
726     this epochs was determined by the minimal inter-stimulus interval beforehand until the minimal
727     inter-trial interval afterwards. Data was re-referenced to the grand average, which allowed us to
728     recover the reference as channel FCz, and then band-pass filtered using a two-pass 4th order
729     Butterworth IIR filter (Fieldtrip default) in the range of 0.5–35 Hz. These filter settings allowed us to
730     distinguish the delta, theta, alpha, and beta band, while filtering out residual high-frequency MR
731     noise. This low-pass filter cut-off was different from a previous analysis of this data in which we set it
732     at 15 Hz (32) because in this analysis, we had a hypothesis on outcome valence encoding in the beta
733     range. We then applied linear baseline correction based on the 200 ms prior to cue onset and used
734     ICA to detect and reject independent components related to eye-blinks, saccades, head motion, and
735     residual MR artifacts (mean number of rejected components per participant: 32.694, range 24–45).
736     Afterwards, we manually rejected trials with residual motion (for all 36 participants: $M$ = 117.722,
737     range 11–499). Based on trial rejection, four participants for which more than 211 (33%) of trials

738    were rejected were excluded from any further analyses (rejected trials after excluding those
739    participants: $M = 81.875$, range 11–194). Finally, we computed a Laplacian filter with the spherical
740    spline method to remove global noise (using the exact electrode positions recorded with Polhemus
741    FASTRAK), which we also used to interpolate previously rejected channels. This filter attenuates
742    more global signals (e.g., signal from deep sources or global noise) and noise (heart-beat and muscle
743    artifacts) while accentuating more local effects (e.g., superficial sources).

## EEG TF decomposition

745    We decomposed the trial-by-trial EEG time series into their time-frequency representations using
746    33 Hanning tapers between 1 and 33 Hz in steps of 1 Hz, every 25 ms from -1000 until 1,300 ms
747    relative to outcome onset. We first zero-padded trials to a length of 8 sec. and then performed time-
748    frequency decomposition in steps of 1 Hz by multiplying the Fourier transform of the trail with the
749    Fourier transform of a Hanning taper of 400 ms width, centered around the time point of interest.
750    This procedure results in an effective resolution of 2.5 Hz (Rayleigh frequency), interpolated in 1 Hz
751    steps, which is more robust to the choice of exact frequency bins. To exclude the possibility of slow
752    drifts in power over the time course of the experiment, we performed baseline correction across
753    participants and trials by fitting a linear model for each channel/ frequency combination with trial
754    number as predictor and the average power 250–50 ms before outcome onset as outcome, and
755    subtracting the power predicted by this model from the data. This procedure is able to remove slow
756    linear drifts in power over time from the data. In absence of such drifts, it is equivalent to correcting
757    all trials by the grand mean across trials per frequency in the selected baseline time window.
758    Afterwards, we averaged power over trials within each condition spanned by performed action (Go/
759    NoGo) and outcome (reward/ no reward/ no punishment/ punishment). We finally converted the
760    average time-frequency data per condition to decibel to ensure that data across frequencies, time
761    points, electrodes, and participants were on same scale.

## EEG analyses

763    All analyses were performed on the average signal of a-priori selected channels Fz, FCz, and
764    Cz based on (9, 32). We again performed model-free and model-based analyses. For the model-free
765    analyses, we sorted trials based on the performed action (Go/ NoGo) and obtained outcome
766    (reward/ no reward/ no punishment/ punishment) and computed the mean TF power across trials
767    for each of the resultant eight conditions for each participant. We tested whether theta power
768    (average power 4–8 Hz) and beta power (average power 13–30 Hz) encoded outcome valence by
769    contrasting favorable (reward/ no punishment) and unfavorable (no reward/ punishment) conditions
770    (irrespective of the performed action). We also tested for differences between Go and NoGo
771    responses in the lower alpha band (6–10 Hz). For all contrasts, we employed two-sided cluster-based
772    permutation tests in a window from 0–1,000 ms relative to outcome onset. For beta power, results
773    were driven by a cluster that was at the edge of 1,000 ms; to more accurately report the time span
774    during which this cluster exceeded the threshold, we extended the time window to 1,300 ms in this
775    particular analysis. Such tests are able to reject the null hypothesis of exchangeability of two
776    experimental conditions, but they are not suited to precisely locate clusters in time-frequency space.
777    Hence, interpretations are mostly based on the visual inspection of plots of the signal time courses.
778    For model-based analyses, similar to fMRI analyses, we used the group-level parameters
779    from the best fitting computational model M5 to compute the trial-by-trial biased PE term and
780    decomposed it into the standard PE term and the difference to the biased PE term. We used both
781    terms as predictors in a multiple linear regression for each channel-time-frequency bin for each

782    participant, and then performed one-sample cluster-based permutation-tests across the resultant $b$-
783    maps of all participants (*78*). For further details on this procedure, see fMRI-inspired EEG analyses.

## fMRI-informed EEG analyses

785    The BOLD signal is sluggish. It is thus hard to determine when different brain regions become
786    active. In contrast, EEG provides much higher temporal resolution. A fruitful approach can be to
787    identify distinct EEG correlates of the BOLD signal in different regions, allowing to test hypotheses
788    about the temporal order in which regions might become active and modulated EEG power (*32, 63*).
789    Furthermore, by using the BOLD signal from different regions in a multiple linear regression, one can
790    control for variance shared among regions (e.g., changes in global signal; variance due to task
791    regressors) and test which region is the best unique predictor of a certain EEG signal. In such an
792    analysis, any correlation between EEG and BOLD signal from a certain region reflects an association
793    above and beyond those induced by task conditions.
794    We used the trial-by-trial BOLD signal in selected regions in a multiple linear regression to predict
795    EEG signal over the scalp (*32, 63*) (building on existing code from https://github.com/tuhauser/TAfT).
796    As a first step, we extracted the volume-by-volume signal (first eigenvariate) from each of the seven
797    regions identified to encode biased PEs (conjunction of $PE_{STD}$ and $PE_{DIF}$: striatum, ACC, vmPFC, left
798    motor cortex, PCC, left ITG, and primary visual cortex). We applied a highpass-filter at 128 s and
799    regressed out nuisance regressors (6 realignment parameters, CSF, OOB, single volumes with strong
800    motion, same as in the fMRI GLM). We then upsampled the signal by a factor 10, epoched it into
801    trials of 8 s duration, and fitted a separate HRF (based on the SPM template) to each trial (58
802    upsampled data points), resulting in trial-by-trial regression weights reflecting the respective BOLD
803    response. We then combined the regression weights of all trials and regions of a certain participant
804    into a design matrix with trials as rows and the seven ROIs as columns, which we used to predict
805    power at each time-frequency-channel bin. As further control variables, we added the behavioral
806    $PE_{STD}$ and $PE_{DIF}$ regressors to the design matrix. All predictors and outcomes were demeaned such
807    that the intercept became zero. Such a multiple linear regression was performed for each participant,
808    resulting in a time-frequency-channel-ROI $b$-map reflecting the association between trial-by-trial
809    BOLD signal and TF power at each time-frequency-channel bin. $B$-maps were Fisher-$z$ transformed,
810    which makes the sampling distribution of correlation coefficients approximately normal and allows
811    for combining them across participants, and analyzed with a cluster-based one-sample permutation
812    $t$-test (*78*) on the mean regression weights over channels Fz, FCz, and Cz across participants in the
813    range of 0–1000 ms, 1–33 Hz. We first obtained a null distribution of maximal cluster mass statistics
814    from 10000 permutations. For each permutation, we flipped the sign of the $b$-map of a random
815    subset of participants, computed a separate $t$-test at each time-frequency bin (bins of 25 ms, 1 Hz)
816    across participants (results in $t$-map), thresholded these maps at $|t| > 2$, and finally computed the
817    maximal cluster mask statistic (sum of all $t$-values) for any cluster (adjacent voxels above threshold).
818    Afterwards, we computed the same $t$-map for the real data, identified the cluster with the biggest
819    cluster-mass statistic, and computed the corresponding $p$-value as number of permutations in the
820    null distribution that were larger than the maximal cluster mass statistic in the real data.

## EEG-informed fMRI analyses

822    For the EEG-informed fMRI analyses, we fit three additional GLMs for which we entered the
823    trial-by-trial theta/ delta power (1–8 Hz), beta power (13–30 Hz), and lower alpha band power (6–10
824    Hz) as parametric regressors on top of the task regressors of the model-free GLM. These measures
825    were created by using the 3-D (time-frequency-channel) $t$-map obtained when contrasting positive

826   vs. negative outcomes (theta/ delta and beta) and Go vs. NoGo conditions (lower alpha band) as a
827   linear filter. We enforced strict frequency cut-offs. For lower alpha band and beta, we used
828   midfrontal channels (Fz/ FCz/ Cz). For theta/ delta power, given the topography that reached far
829   beyond midfrontal channels and over the entire frontal scalp, we used a much wider ROI (AF3/ AF4/
830   AF7/ AF8/ F1/ F2/ F3/ F4/ F5/ F6/ F7/ F8/ FC1/ FC2/ FC3/ FC4/ FC5/ FC6/ FCz/ Fp1/ Fp2/ Fpz/ Fz). We
831   extracted those maps and retained all voxels with t > 2. These masks were applied to the trial-by-trial
832   time-frequency data to create weighted summary measures of the average power in the identified
833   clusters in each trial. For trials for which EEG data was rejected, we imputed the participant mean
834   value of the respective action (Go/ NoGo) x outcome (reward/ no reward/ no punishment/
835   punishment) condition. Note that this approach accentuates differences between conditions, which
836   are already captured by the task regressors in the GLM, but decreases trial-by-trial variability within
837   each condition, which is of interest in this analysis. This imputation approach is thus conservative.
838   While trial-by-trial beta and theta power were largely uncorrelated, mean $r$ = 0.104, range -0.118–
839   0.283 across participants, and so were beta and alpha, mean $r$ = 0.097, range -0.162–0.284 across
840   participants, theta and alpha power moderately correlate, mean $r$ = 0.412, range 0.121–0.836 across
841   participants, warranting the use of a separate channel ROI for theta and using separate GLMs for
842   each frequency band.

### Analyses of behavior as a function of BOLD signal and EEG power

844   We used mixed-effects logistic regression to analyze "stay behavior", i.e., whether
845   participants repeated an action on the next encounter of the same cue, as a function of BOLD signal
846   and EEG power in selected regions. For analyses featuring BOLD signal, we used the trial-by-trial HRF
847   amplitude also used for fMRI-informed EEG analyses. For analyses featuring EEG, we used the trial-
848   by-trial EEG power also used in the EEG-informed fMRI analyses.

## References

850   1.    P. Dayan, Y. Niv, B. Seymour, N. Daw, The misbehavior of value and the discipline of the will.
851         *Neural Networks*. **19**, 1153–1160 (2006).

852   2.    M. Guitart-Masip, E. Duzel, R. Dolan, P. Dayan, Action versus valence in decision making.
853         *Trends Cogn. Sci.* **18**, 194–202 (2014).

854   3.    J. C. Swart, M. I. Froböse, J. L. Cook, D. E. M. Geurts, M. J. Frank, R. Cools, H. E. den Ouden,
855         Catecholaminergic challenge uncovers distinct Pavlovian and instrumental mechanisms of
856         motivated (in)action. *Elife*. **6**, 1–54 (2017).

857   4.    A. Mkrtchian, J. Aylward, P. Dayan, J. P. Roiser, O. J. Robinson, Modeling avoidance in mood
858         and anxiety disorders using reinforcement learning. *Biol. Psychiatry*. **82**, 532–539 (2017).

859   5.    Q. J. M. Huys, M. Gölzer, E. Friedel, A. Heinz, R. Cools, P. Dayan, R. J. Dolan, The specificity of
860         Pavlovian regulation is associated with recovery from depression. *Psychol. Med.* **46**, 1027–
861         1035 (2016).

862   6.    Q. J. M. Huys, R. Cools, M. Gölzer, E. Friedel, A. Heinz, R. J. Dolan, P. Dayan, Disentangling the
863         roles of approach, activation and valence in instrumental and Pavlovian responding. *PLoS
864         Comput. Biol.* **7**, e1002028 (2011).

865   7.    Y.-L. Boureau, P. Sokol-Hessner, N. D. Daw, Deciding how to decide: Self-control and meta-
866         decision making. *Trends Cogn. Sci.* **19**, 700–710 (2015).

867   8.    M. Guitart-Masip, Q. J. M. Huys, L. Fuentemilla, P. Dayan, E. Duzel, R. J. Dolan, Go and no-go

868    learning in reward and punishment: Interactions between affect and effect. *Neuroimage*. **62**,
869    154–166 (2012).

870  9.   J. C. Swart, M. J. Frank, J. I. Määttä, O. Jensen, R. Cools, H. E. M. den Ouden, Frontal network
871       dynamics reflect neurocomputational mechanisms for reducing maladaptive biases in
872       motivated action. *PLOS Biol.* **16**, e2005979 (2018).

873  10.  L. de Boer, J. Axelsson, R. Chowdhury, K. Riklund, R. J. Dolan, L. Nyberg, L. Bäckman, M.
874       Guitart-Masip, Dorsal striatal dopamine D1 receptor availability predicts an instrumental bias
875       in action learning. *Proc. Natl. Acad. Sci.* **116**, 261–270 (2019).

876  11.  D. R. Williams, H. Williams, Auto-maintenance in the pigeon: Sustained pecking despite
877       contingent non-reinforcement. *J. Exp. Anal. Behav.* **12**, 511–520 (1969).

878  12.  P. L. Brown, H. M. Jenkins, Autoshaping of pigeon's key-peck. *J. Exp. Anal. Behav.* **11**, 1–8
879       (1968).

880  13.  I. Ritov, J. Baron, Reluctance to vaccinate: Omission bias and ambiguity. *J. Behav. Decis. Mak.*
881       **3**, 263–277 (1990).

882  14.  M. Zeelenberg, J. van der Pligt, N. K. de Vries, Attributions of responsibility and affective
883       reactions to decision outcomes. *Acta Psychol. (Amst).* **104**, 303–315 (2000).

884  15.  M. J. Frank, Dynamic dopamine modulation in the basal ganglia: A neurocomputational
885       account of cognitive deficits in medicated and nonmedicated Parkinsonism. *J. Cogn. Neurosci.*
886       **17**, 51–72 (2005).

887  16.  A. G. E. Collins, M. J. Frank, Opponent actor learning (OpAL): Modeling interactive effects of
888       striatal dopamine on reinforcement learning and choice incentive. *Psychol. Rev.* **121**, 337–366
889       (2014).

890  17.  B. B. Doll, W. J. Jacobs, A. G. Sanfey, M. J. Frank, Instructional control of reinforcement
891       learning: A behavioral and neurocomputational investigation. *Brain Res.* **1299**, 74–94 (2009).

892  18.  L. Y. Atlas, B. B. Doll, J. Li, N. D. Daw, E. A. Phelps, Instructed knowledge shapes feedback-
893       driven aversive learning in striatum and orbitofrontal cortex, but not the amygdala. *Elife.* **5**, 1–
894       26 (2016).

895  19.  N. D. Daw, S. J. Gershman, B. Seymour, P. Dayan, R. J. Dolan, Model-based influences on
896       humans' choices and striatal prediction errors. *Neuron.* **69**, 1204–1215 (2011).

897  20.  S. W. Lee, S. Shimojo, J. P. O'Doherty, Neural computations underlying arbitration between
898       model-based and model-free learning. *Neuron.* **81**, 687–699 (2014).

899  21.  P. Piray, I. Toni, R. Cools, Human choice strategy varies with anatomical projections from
900       ventromedial prefrontal cortex to medial striatum. *J. Neurosci.* **36**, 2857–2867 (2016).

901  22.  T. E. J. Behrens, M. W. Woolrich, M. E. Walton, M. F. S. Rushworth, Learning the value of
902       information in an uncertain world. *Nat. Neurosci.* **10**, 1214–1221 (2007).

903  23.  D. Meder, N. Kolling, L. Verhagen, M. K. Wittmann, J. Scholl, K. H. Madsen, O. J. Hulme, T. E. J.
904       Behrens, M. F. S. Rushworth, Simultaneous representation of a spectrum of dynamically
905       changing value estimates during decision making. *Nat. Commun.* **8**, 1942 (2017).

906  24.  A. J. van Nuland, R. C. Helmich, M. F. Dirkx, H. Zach, I. Toni, R. Cools, H. E. M. den Ouden,
907       Effects of dopamine on reinforcement learning in Parkinson's disease depend on motor
908       phenotype. *Brain*, 1–13 (2020).

909  25.  I. van de Vijver, K. R. Ridderinkhof, M. X. Cohen, Frontal oscillatory dynamics predict feedback

910         learning and action adjustment. *J. Cogn. Neurosci.* **23**, 4106–4121 (2011).

911   26.   M. X. Cohen, K. A. Wilmes, I. van de Vijver, Cortical electrophysiological network dynamics of
912         feedback learning. *Trends Cogn. Sci.* **15**, 558–566 (2011).

913   27.   J. F. Cavanagh, M. J. Frank, T. J. Klein, J. J. B. Allen, Frontal theta links prediction errors to
914         behavioral adaptation in reinforcement learning. *Neuroimage.* **49**, 3198–3209 (2010).

915   28.   J. F. Cavanagh, Cortical delta activity reflects reward prediction error and related behavioral
916         adjustments, but at different times. *Neuroimage.* **110**, 205–216 (2015).

917   29.   D. Talmi, R. Atkinson, W. El-Deredy, The feedback-related negativity signals salience
918         prediction errors, not reward prediction errors. *J. Neurosci.* **33**, 8264–8269 (2013).

919   30.   E. M. Bernat, L. D. Nelson, A. R. Baskin-Sommers, Time-frequency theta and delta measures
920         index separable components of feedback processing in a gambling task. *Psychophysiology.* **52**,
921         626–637 (2015).

922   31.   J. Marco-Pallarés, T. F. Münte, A. Rodríguez-Fornells, The role of high-frequency oscillatory
923         activity in reward processing and learning. *Neurosci. Biobehav. Rev.* **49**, 1–7 (2015).

924   32.   J. Algermissen, J. C. Swart, R. Scheeringa, R. Cools, H. E. M. den Ouden, Striatal BOLD signal
925         and midfrontal theta oscillations express motivation for action. *Cereb. Cortex* (2021).

926   33.   J. Cockburn, A. G. E. Collins, M. J. Frank, A reinforcement learning mechanism responsible for
927         the valuation of free choice. *Neuron.* **83**, 551–557 (2014).

928   34.   B. C. Wittmann, N. D. Daw, B. Seymour, R. J. Dolan, Striatal activity underlies novelty-based
929         choice in humans. *Neuron.* **58**, 967–973 (2008).

930   35.   O. Bartra, J. T. McGuire, J. W. Kable, The valuation system: A coordinate-based meta-analysis
931         of BOLD fMRI experiments examining neural correlates of subjective value. *Neuroimage.* **76**,
932         412–427 (2013).

933   36.   M. J. Frank, B. S. Woroch, T. Curran, Error-related negativity predicts reinforcement learning
934         and conflict biases. *Neuron.* **47**, 495–501 (2005).

935   37.   M. J. Sharpe, C. Y. Chang, M. A. Liu, H. M. Batchelor, L. E. Mueller, J. L. Jones, Y. Niv, G.
936         Schoenbaum, Dopamine transients are sufficient and necessary for acquisition of model-
937         based associations. *Nat. Neurosci.* **20**, 735–742 (2017).

938   38.   M. J. Sharpe, T. Stalnaker, N. W. Schuck, S. Killcross, G. Schoenbaum, Y. Niv, An integrated
939         model of action selection: Distinct modes of cortical control of striatal decision making. *Annu.*
940         *Rev. Psychol.* **70**, 53–76 (2019).

941   39.   A. Pasupathy, E. K. Miller, Different time courses of learning-related activity in the prefrontal
942         cortex and striatum. *Nature.* **433**, 873–876 (2005).

943   40.   E. G. Antzoulatos, E. K. Miller, Increases in functional connectivity between prefrontal cortex
944         and striatum during category learning. *Neuron.* **83**, 216–225 (2014).

945   41.   M. Seo, E. Lee, B. B. Averbeck, Action selection and action value in frontal-striatal circuits.
946         *Neuron.* **74**, 947–960 (2012).

947   42.   J. D. Howard, R. Reynolds, D. E. Smith, J. L. Voss, G. Schoenbaum, T. Kahnt, Targeted
948         stimulation of human orbitofrontal networks disrupts outcome-guided behavior. *Curr. Biol.*
949         **30**, 490-498.e4 (2020).

950   43.   M. R. van Schouwenburg, J. O'Shea, R. B. Mars, M. F. S. Rushworth, R. Cools, Controlling

951      human striatal cognitive function via the frontal cortex. *J. Neurosci.* **32**, 5631–5637 (2012).

952   44.   E. Fouragnan, C. Retzler, K. Mullinger, M. G. Philiastides, Two spatiotemporally distinct value
953      systems shape reward-based learning in the human brain. *Nat. Commun.* **6**, 8107 (2015).

954   45.   W. H. Alexander, J. W. Brown, Medial prefrontal cortex as an action-outcome predictor. *Nat.*
955      *Neurosci.* **14**, 1338–1344 (2011).

956   46.   W. H. Alexander, J. W. Brown, Frontal cortex function as derived from hierarchical predictive
957      coding. *Sci. Rep.* **8**, 3843 (2018).

958   47.   P. Enel, J. D. Wallis, E. L. Rich, Stable and dynamic representations of value in the prefrontal
959      cortex. *Elife.* **9**, 1–23 (2020).

960   48.   S. Vyas, D. J. O'Shea, S. I. Ryu, K. V. Shenoy, Causal role of motor preparation during error-
961      driven learning. *Neuron*, 1–11 (2020).

962   49.   R. Shadmehr, M. A. Smith, J. W. Krakauer, Error correction, sensory prediction, and adaptation
963      in motor control. *Annu. Rev. Neurosci.* **33**, 89–108 (2010).

964   50.   S. Sadaghiani, R. Scheeringa, K. Lehongre, B. Morillon, A.-L. Giraud, A. Kleinschmidt, Intrinsic
965      connectivity networks, alpha oscillations, and tonic alertness: A simultaneous
966      electroencephalography/functional magnetic resonance imaging study. *J. Neurosci.* **30**,
967      10243–10250 (2010).

968   51.   C. Andreou, H. Frielinghaus, J. Rauh, M. Mußmann, S. Vauth, P. Braun, G. Leicht, C. Mulert,
969      Theta and high-beta networks for feedback processing: A simultaneous EEG-fMRI study in
970      healthy male subjects. *Transl. Psychiatry.* **7**, e1016–e1016 (2017).

971   52.   J. Marco-Pallarés, D. Cucurell, T. Cunillera, R. García, A. Andrés-Pueyo, T. F. Münte, A.
972      Rodríguez-Fornells, Human oscillatory activity associated to reward processing in a gambling
973      task. *Neuropsychologia.* **46**, 241–248 (2008).

974   53.   R. Scheeringa, M. C. M. Bastiaansen, K. M. Petersson, R. Oostenveld, D. G. Norris, P. Hagoort,
975      Frontal theta EEG activity correlates negatively with the default mode network in resting
976      state. *Int. J. Psychophysiol.* **67**, 242–251 (2008).

977   54.   R. Scheeringa, K. M. Petersson, R. Oostenveld, D. G. Norris, P. Hagoort, M. C. M. Bastiaansen,
978      Trial-by-trial coupling between EEG and BOLD identifies networks related to alpha and theta
979      EEG power increases during working memory maintenance. *Neuroimage.* **44**, 1224–1238
980      (2009).

981   55.   E. Fouragnan, C. Retzler, M. G. Philiastides, Separate neural representations of prediction
982      error valence and surprise: Evidence from an fMRI meta-analysis. *Hum. Brain Mapp.* (2018),
983      doi:10.1002/hbm.24047.

984   56.   S. N. Haber, The primate basal ganglia: Parallel and integrative networks. *J. Chem. Neuroanat.*
985      **26**, 317–330 (2003).

986   57.   A. K. Engel, P. Fries, Beta-band oscillations—signalling the status quo? *Curr. Opin. Neurobiol.*
987      **20**, 156–165 (2010).

988   58.   J. Feingold, D. J. Gibson, B. Depasquale, A. M. Graybiel, Bursts of beta oscillation differentiate
989      postperformance activity in the striatum and motor cortex of monkeys performing movement
990      tasks. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 13687–13692 (2015).

991   59.   T. U. Hauser, R. Iannaccone, P. Stämpfli, R. Drechsler, D. Brandeis, S. Walitza, S. Brem, The
992      feedback-related negativity (FRN) revisited: New insights into the localization, meaning and

993          network organization. *Neuroimage*. **84**, 159–168 (2014).

994    60.   J. R. Wessel, A. R. Aron, On the globality of motor suppression: Unexpected events and their
995          influence on behavior and cognition. *Neuron*. **93**, 259–280 (2017).

996    61.   N. Trudel, J. Scholl, M. C. Klein-Flügge, E. Fouragnan, L. Tankelevitch, M. K. Wittmann, M. F. S.
997          Rushworth, Polarity of uncertainty representation during exploration and exploitation in
998          ventromedial prefrontal cortex. *Nat. Hum. Behav.* (2020), doi:10.1038/s41562-020-0929-3.

999    62.   P. Domenech, S. Rheims, E. Koechlin, Neural mechanisms resolving exploitation-exploration
1000         dilemmas in the medial prefrontal cortex. *Science (80-. ).* **369**, eabb0184 (2020).

1001   63.   T. U. Hauser, L. T. Hunt, R. Iannaccone, S. Walitza, D. Brandeis, S. Brem, R. J. Dolan,
1002         Temporally dissociable contributions of human medial prefrontal subregions to reward-
1003         guided learning. *J. Neurosci.* **35**, 11209–11220 (2015).

1004   64.   D. Foti, A. Weinberg, J. Dien, G. Hajcak, Event-related potential activity in the basal ganglia
1005         differentiates rewards from nonrewards: Temporospatial principal components analysis and
1006         source localization of the feedback negativity. *Hum. Brain Mapp.* **32**, 2207–2216 (2011).

1007   65.   M. X. Cohen, J. F. Cavanagh, H. A. Slagter, Event-related potential activity in the basal ganglia
1008         differentiates rewards from nonrewards: Temporospatial principal components analysis and
1009         source localization of the feedback negativity: Commentary. *Hum. Brain Mapp.* **32**, 2270–2271
1010         (2011).

1011   66.   K. Amemori, S. Amemori, D. J. Gibson, A. M. Graybiel, Striatal microstimulation induces
1012         persistent and repetitive negative decision-making predicted by striatal beta-band oscillation.
1013         *Neuron*. **99**, 829-841.e6 (2018).

1014   67.   K. Amemori, S. Amemori, D. J. Gibson, A. M. Graybiel, Striatal beta oscillation and neuronal
1015         activity in the primate caudate nucleus differentially represent valence and arousal under
1016         approach-avoidance conflict. *Front. Neurosci.* **14**, 1–17 (2020).

1017   68.   J. F. Cavanagh, I. Eisenberg, M. Guitart-Masip, Q. J. M. Huys, M. J. Frank, Frontal theta
1018         overrides Pavlovian learning biases. *J. Neurosci.* **33**, 8541–8548 (2013).

1019   69.   D. J. Barr, R. Levy, C. Scheepers, H. J. Tily, Random effects structure for confirmatory
1020         hypothesis testing: Keep it maximal. *J. Mem. Lang.* **68**, 255–278 (2013).

1021   70.   R. A. Rescorla, A. R. Wagner, in *Classical Conditioning II⬚:Current Research and Theory*, A. H.
1022         Black, W. F. Prokasy, Eds. (Appleton-Century-Crofts., New York, NY, 1972;
1023         http://homepage.mac.com/sanagnos/rescorlawagner1972.pdf), vol. 21, pp. 64–99.

1024   71.   P. Piray, A. Dezfouli, T. Heskes, M. J. Frank, N. D. Daw, Hierarchical Bayesian inference for
1025         concurrent model fitting and comparison for group studies. *PLOS Comput. Biol.* **15**, e1007043
1026         (2019).

1027   72.   K. E. Stephan, W. D. Penny, J. Daunizeau, R. J. Moran, K. J. Friston, Bayesian model selection
1028         for group studies. *Neuroimage*. **46**, 1004–1017 (2009).

1029   73.   F. Krause, C. Benjamins, J. Eck, M. Lührs, R. Hoof, R. Goebel, Active head motion reduction in
1030         magnetic resonance imaging using tactile feedback. *Hum. Brain Mapp.* **40**, 4026–4037 (2019).

1031   74.   R. H. R. Pruim, M. Mennes, D. van Rooij, A. Llera, J. K. Buitelaar, C. F. Beckmann, ICA-AROMA:
1032         A robust ICA-based strategy for removing motion artifacts from fMRI data. *Neuroimage*. **112**,
1033         267–277 (2015).

1034   75.   R. C. Wilson, Y. Niv, Is model fitting necessary for model-based fMRI? *PLOS Comput. Biol.* **11**,

1035        e1004237 (2015).

1036   76.  P. J. Allen, O. Josephs, R. Turner, A method for removing imaging artifact from continuous EEG
1037        recorded during functional MRI. *Neuroimage.* **12**, 230–239 (2000).

1038   77.  R. Oostenveld, P. Fries, E. Maris, J.-M. Schoffelen, FieldTrip: Open source software for
1039        advanced analysis of MEG, EEG, and invasive electrophysiological data. *Comput. Intell.*
1040        *Neurosci.* **2011**, 1–9 (2011).

1041   78.  L. T. Hunt, M. W. Woolrich, M. F. S. Rushworth, T. E. J. Behrens, Trial-type dependent frames
1042        of reference for value comparison. *PLoS Comput. Biol.* **9**, e1003225 (2013).

1043   79.  S. Palminteri, V. Wyart, E. Koechlin, The importance of falsification in computational cognitive
1044        modeling. *Trends Cogn. Sci.* **21**, 425–433 (2017).

1045   80.  M. R. Nassar, M. J. Frank, Taming the beast: Extracting generalizable knowledge from
1046        computational models of cognition. *Curr. Opin. Behav. Sci.* **11**, 49–54 (2016).

1047   81.  M. X. Cohen, T. H. Donner, Midfrontal conflict-related theta-band power reflects neural
1048        oscillations that predict behavior. *J. Neurophysiol.* **110**, 2752–2763 (2013).

1049   82.  G. H. Proudfit, The reward positivity: From basic research on reward to a biomarker for
1050        depression. *Psychophysiology.* **52**, 449–459 (2015).

1051   83.  K. Paul, E. Vassena, M. C. Severo, G. Pourtois, Dissociable effects of reward magnitude on
1052        fronto-medial theta and FRN during performance monitoring. *Psychophysiology* (2019),
1053        doi:10.1111/psyp.13481.

1054   84.  T. D. Sambrook, J. Goslin, Principal components analysis of reward prediction errors in a
1055        reinforcement learning task. *Neuroimage.* **124**, 276–286 (2016).

1056   85.  N. Yeung, A. G. Sanfey, Independent coding of reward magnitude and valence in the human
1057        brain. *J. Neurosci.* **24**, 6258–6264 (2004).

1058   86.  L. Kreussel, J. Hewig, N. Kretschmer, H. Hecht, M. G. H. Coles, W. H. R. Miltner, The influence
1059        of the magnitude, probability, and valence of potential wins and losses on the amplitude of
1060        the feedback negativity. *Psychophysiology.* **49**, 207–219 (2012).

1061   87.  A. Sato, A. Yasuda, H. Ohira, K. Miyawaki, M. Nishikawa, H. Kumano, T. Kuboki, Effects of value
1062        and reward magnitude on feedback negativity and P300. *Neuroreport.* **16**, 407–411 (2005).

1063   88.  D. Tanner, K. Morgan-Short, S. J. Luck, How inappropriate high-pass filters can produce
1064        artifactual effects and incorrect conclusions in ERP studies of language and cognition.
1065        *Psychophysiology.* **52**, 997–1009 (2015).

1066   89.  Y. Wu, X. Zhou, The P300 and reward valence, magnitude, and expectancy in outcome
1067        evaluation. *Brain Res.* **1286**, 114–122 (2009).

1068   90.  R. Scheeringa, P. Fries, K.-M. Petersson, R. Oostenveld, I. Grothe, D. G. Norris, P. Hagoort, M.
1069        C. M. Bastiaansen, Neuronal dynamics underlying high-and low-frequency EEG oscillations
1070        contribute independently to the human BOLD signal. *Neuron.* **69**, 572–583 (2011).

1071   91.  J. M. Zumer, R. Scheeringa, J.-M. Schoffelen, D. G. Norris, O. Jensen, Occipital alpha activity
1072        during stimulus processing gates the information flow to object-selective cortex. *PLoS Biol.* **12**,
1073        e1001965 (2014).

1074   92.  M. T. Jurkiewicz, W. C. Gaetz, A. C. Bostan, D. Cheyne, Post-movement beta rebound is
1075        generated in motor cortex: Evidence from neuromagnetic recordings. *Neuroimage.* **32**, 1281–
1076        1289 (2006).

1077    93.    P. Ritter, M. Moosmann, A. Villringer, Rolandic alpha and beta EEG rhythms' strengths are
1078           inversely related to fMRI-BOLD signal in primary somatosensory and motor cortex. *Hum. Brain*
1079           *Mapp.* **30**, 1168–1187 (2009).

1080    94.    W. Klimesch, EEG alpha and theta oscillations reflect cognitive and memory performance: A
1081           review and analysis. *Brain Res. Rev.* **29**, 169–195 (1999).

## Author contributions

Conceptualization: JA, JCS, RC, HEMDO
Data curation: JA
Formal analysis: JA
Funding acquisition: JCS, RC, HEMDO
Investigation: JA, JCS
Methodology: JA, HEMDO
Project administration: JA, JCS, HEMDO
Resources: RC, HEMDO
Software: JA, JCS, HEMDO
Supervision: JCS, RS, RC, HEMDO
Validation: JA, JCS, RS, RC, HEMDO
Visualization: JA
Writing – original draft: JA, HEMDO
Writing – review & editing: JA, JCS, RS, RC, HEMDO

## Competing interests

Authors declare that they have no competing interests.

## Data availability statement

All data and code will be made available upon manuscript acceptance.
Group-level unthresholded fMRI z-maps are available on Neurovault
(https://neurovault.org/collections/11184/).

1116 # Supplementary Materials to "Biased

1117 # credit assignment in motivational

1118 # learning biases arises through prefrontal

1119 # influences on striatal learning"

1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150

## S01: Behavioral, fMRI, and EEG analyses with only the 29 participants included in EEG-fMRI analyses

We repeated the behavioral, fMRI, and EEG analyses reported in the main text while excluding the seven participants that were also not included in the fMRI-inspired EEG analyses in the main text: (a) two participants due to fMRI co-registration failure, which were also not included in the fMRI-only analyses; (b) four further participants who exhibited excessive residual noise in their EEG data (> 33% rejected trials) and were thus also not included in the EEG-only analyses, and finally (c) one more participant who (together with four other participants already excluded) exhibited regression weights for every regressor about ten times larger than for other participants.

Participants in this subgroup learned the task, reflected in a significant main effect of required action on responses, $b = 0.896$, $SE = 0.129$, $\chi^2(1) = 28.398$, $p < .001$, and exhibited motivational biases, reflected in a significant main effect of cue valence on responses, $b = 0.439$, $SE = 0.084$, $\chi^2(1) = 19.308$, $p < .001$. The interaction between required action and cue valence was not significant, $b = 0.025$, $SE = 0.085$, $\chi^2(1) = 0.111$, $p = .739$.

Participants in this subgroup also showed biased learning: They were more likely to repeat an action after a favorable outcome (main effect of outcome valence: $b = .0553$, $SE = 0.059$, $\chi^2(1) = 40.920$, $p < .001$. After salient outcomes, they adjusted their responses more strongly after feedback on Go than on NoGo responses, in line with our model of biased learning and as reflected in a significant three-way interaction between action, salience, and valence, $b = 0.266$, $SE = 0.055$, $\chi^2(1) = 16.862$, $p < .001$. When only analyzing trials with salient outcomes, outcome valence was more likely to affect response repetition following Go relative to NoGo responses, $b = 0.324$, $SE = 0.079$, $\chi^2(1) = 13.266$, $p < .001$, with a stronger effect of outcome valence after Go responses, $b = 1.342$, $SE = 0.120$, $\chi^2(1) = 49.003$, $p = .001$, than NoGo responses, $b = 0.693$, $SE = 0.129$, $\chi^2(1) = 18.988$, $p < .001$.

In this subgroup of participants, Bayesian model selection clearly favored the full asymmetric pathways models featuring response and learning biases (M5, model frequency: 81.81%, protected exceedance probability: 100%). In sum, behavioral results were qualitatively identical when analyzing only this subgroup of only 29 participants.
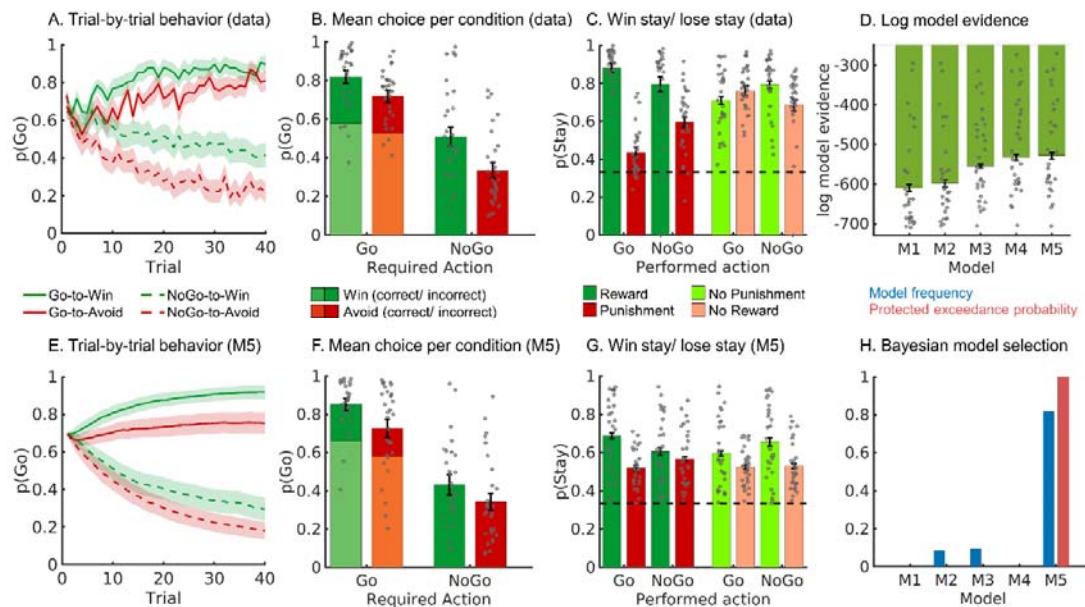
**Figure S01A. Behavioral performance in the subgroup of 29 participants included in the fMRI-inspired EEG analyses.** (A) Trial-by-trial proportion of Go responses (±SEM across participants) for Go cues (solid lines) and NoGo cues (dashed lines). The motivational bias is already present from very early trials onwards, as participants made more Go responses to Win than Avoid cues (i.e., green lines are above red lines). Additionally, participants clearly learn whether to make a Go response or not (proportion of Go responses increases for Go cues and decreases for NoGo cues). (B) Mean (±SEM across participants) proportion Go responses per cue condition (points are individual participants' means). C) Probability to repeat a response ("stay") on the next encounter of the same cue as a function of action and outcome. Learning is reflected in higher probability of staying after positive outcomes than after negative outcomes (main effect of outcome valence). Biased learning is evident in learning from salient outcomes, where this valence effect was stronger after Go responses than NoGo responses. Dashed line indicates chance level choice ($p_{Stay}$ = 0.33). (D) Log-model evidence favors the asymmetric pathways model (M5 over simpler models (M1-M4). (E-G) Trial-by-trial proportion of Go responses, mean proportion Go responses, and probability of staying based on one-step-ahead predictions using parameters (hierarchical Bayesian inference) of the winning model (asymmetric pathways model, M5). (H) Model frequency and protected exceedance probability indicate best fit for model M5 (asymmetric pathways model), in line with log model evidence.

1179

1180          Regarding fMRI findings, we first repeated the model-free GLM just contrasting favorable and
1181    non-favorable outcomes. BOLD signal was higher for favorable than non-favorable outcomes in five
1182    clusters, namely in vmPFC, striatum, amygdala, and hippocampus ($z_{max}$ = 5.65, $p$ = 2.24e-25, 6110
1183    voxels, MNI coordinates xyz = [6 30 -12]), left superior lateral occipital cortex ($z_{max}$ = 4.40, $p$ = .00144,
1184    367 voxels, xyz = [-46 -68 46]), right occipital pole ($z_{max}$ = 4.45, $p$ = .00154, 363 voxels, xyz = [12 -92 -
1185    12]), posterior cingulate cortex ($z_{max}$ = 4.36, $p$ = .00181, 353 voxels, xyz = [-2 -48 28]), and left middle
1186    temporal gyrus ($z_{max}$ = 4.63, $p$ = .00548, 289 voxels, xyz = [-60 -10 -16]). The clusters in left slOCC, PCC,
1187    and left MTG emerged anew compared to the original analysis comprising 34 participants. Also,
1188    compared to the original analysis, clusters in left orbitofrontal cortex and left superior frontal gyrus
1189    were merged with the cluster in vmPFC. In sum, all clusters from the original analysis were found
1190    back, plus some additional clusters.
1191          There was also one cluster in right orbitofrontal cortex ($z_{max}$ = 4.37, $p$ = .0209, 217 voxels, xyz
1192    = [30 62 -2]) in which BOLD signal was higher for non-favorable than favorable outcomes. Compared
1193    to the original analysis comprising 34 participants, clusters in precuneous and right superior frontal
1194    gyrus were not significant.

1195    In the model-based GLM featuring regressors for standard PEs and the difference term
1196    towards biased PEs, BOLD signal correlated with standard PEs in ten clusters, namely in vmPFC,
1197    striatum, bilateral amygdala and hippocampus ($z_{max}$ = 6.04, $p$ = .4.78e-44, 8848 voxels, xyz = [12 14 -
1198    6]), left superior frontal gyrus ($z_{max}$ = 5.58, $p$ = 3.5e-10, 1043 voxels, xyz = [-18 34 52]), left occipital
1199    pole and lingual gyrus ($z_{max}$ = 6.23, $p$ = 7.18e-10, 998 voxels, xyz = [10 -92 -10]), posterior cingulate
1200    cortex ($z_{max}$ = 5.12, $p$ = 8.57e-10, 987 voxels, xyz = [4 -36 48]), left inferior temporal gyrus ($z_{max}$ = 5.03,
1201    $p$ = 7.07e-09, 859 voxels, xyz = [-52 -46 -10]), right anterior middle temporal gyrus ($z_{max}$ = 5.32, $p$ =
1202    .000292, 314 voxels, xyz = [62 -4 -16]), right cerebellum ($z_{max}$ = 5.32, $p$ = .002228, 231 voxels, xyz = [44
1203    -72 -40]), left superior lateral occipital cortex ($z_{max}$ = 4.69, $p$ = .00322, 218 voxels, xyz = [-46 -74 -38]),
1204    right caudate ($z_{max}$ = 4.33, $p$ = .00538, 199 voxels, xyz = [20 12 22]), and right middle temporal gyrus
1205    ($z_{max}$ = 4.09, $p$ = .0129, 189 voxels, xyz = [54 -38 -12]). The clusters in left superior lateral occipital
1206    cortex, right caudate, and right posterior middle temporal gyrus emerged anew by splitting from
1207    larger clusters visible in the original analysis based on 34 participants. Vice versa, the cluster in left
1208    middle temporal gyrus reported for the original analysis was merged with a bigger cluster in the
1209    analysis of only 29 participants. The clusters in postcentral gyrus and ACC observed in the original
1210    analysis based on 34 participants were not significant anymore; however, they were still visible at a
1211    level of $z$ > 3.1 uncorrected.
1212    BOLD signal correlated significantly negatively with standard PEs in a single cluster in right
1213    superior frontal gyrus ($z_{max}$ = 5.04, $p$ = .00771, 186 voxels, xyz = [6 26 64]), similar to the respective
1214    cluster reported in the original analysis. In contrast, the clusters in right occipital pole, intracalcarine
1215    cortex, and left inferior lateral occipital cortex were not significant any more, though visible at a level
1216    of $z$ > 3.1 uncorrected.
1217    BOLD signal in six clusters correlated significantly positively with the difference term towards
1218    biased PEs, namely in large parts of cortex and subcortex including striatum ($z_{max}$ = 6-54, $p$ = 0, 29428
1219    voxels, xyz = [34 -84 20]), dorsomedial prefrontal cortex ($z_{max}$ = 5.94, $p$ = 2.69e-40, 7001 voxels, xyz =
1220    [6 22 34]), right insula ($z_{max}$ = 5.76, $p$ = 7.84e-27, 3847 voxels, xyz = [34 20 -8]), thalamus and
1221    brainstem ($z_{max}$ = 5.10, $p$ = 4.06e-18, 2169 voxels, xyz = [4 -30 0]), left caudate ($z_{max}$ = 4.71, $p$ =
1222    .000188, 305 voxels, xyz = [-12 8 6]) and another cluster in brainstem ($z_{max}$ = 4.05, $p$ = .0151, 160
1223    voxels, xyz = [4 -30 -30]). Clusters in dmPFC, right insula, and left caudate split from larger clusters
1224    reported in the original analysis. Vice versa, the cluster in left insula reported in the original analysis
1225    merged with the largest cluster. The clusters in right middle temporal gyrus and right insula were
1226    missing in the analysis of only 29 participants, but visible at a level of $z$ > 3.1 uncorrected.
1227    BOLD signal in three clusters correlated significantly negatively with the difference term
1228    towards biased PEs, namely in vmPFC ($z_{max}$ = 4.23, $p$ = .0051, 185 voxels, xyz = [-12 48 -6]), left
1229    hippocampus ($z_{max}$ = 4.58, $p$ = .00857, 168 voxels, xyz = [-26 -14 -22]), and left medial temporal gyrus
1230    ($z_{max}$ = 4.30, $p$ = .0172, 146 voxels, xyz = [-62 -4 -16]). Compared to the original analysis, the cluster in
1231    vmPFC emerged anew.
1232    When computing the conjunction between both (positive) contrasts, BOLD signal encoded
1233    both the standard and the difference in four clusters, namely in vmPFC, bilateral striatum, bilateral
1234    ITG, and V1. Clusters in ACC, left motor cortex, and PCC were not significant any more (because they
1235    were z > 3.1, but not significant after cluster correction in the standard PE contrast). However, new
1236    (though rather small) clusters of biased PE encoding emerged in right insula, left amygdala, and left
1237    OFC. In sum, results when analyzing only this subgroup of only 29 participants were largely similar to
1238    results based on the full sample; however, clusters of biased PE encoding in left motor cortex, ACC,
1239    and PCC were small and thus did not survive cluster correction in this subgroup.
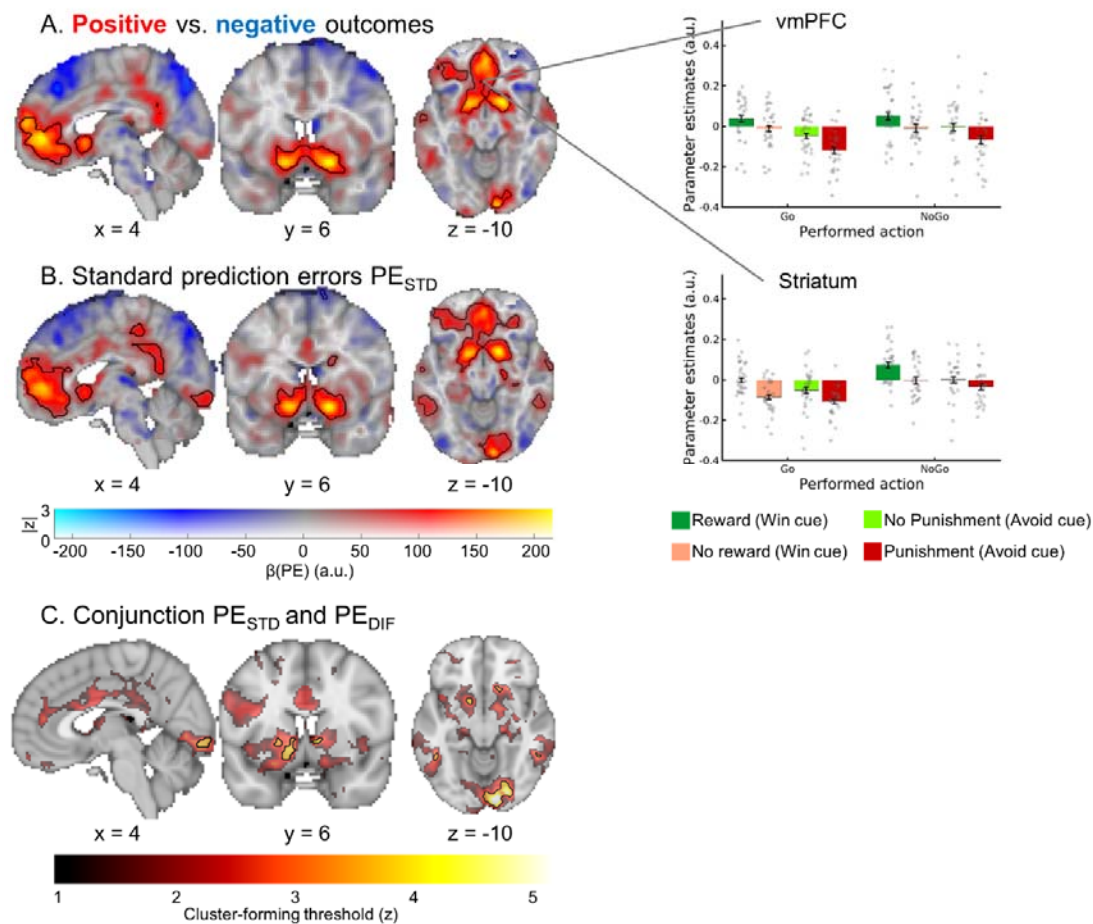
1240



**Figure S01B. BOLD signal reflecting outcome processing in the subgroup of 29 participants included in the fMRI-inspired EEG analyses.** (A) BOLD signal was higher for favorable outcomes (rewards, no punishments) compared with unfavorable outcomes (no rewards, punishments) in a range of regions including bilateral ventral striatum and vmPFC. BOLD effects displayed using a dual-coding data visualization approach with color indicating the parameter estimates and opacity the associated z-statistics. Significant clusters are surrounded by black edges. Bar plots show parameter estimates per action x outcome condition (±SEM across participants) (B) When using the trial-by-trial PEs participants experienced as model-based regressors in our GLM, positive PE correlations occurred in several regions including importantly the ventral striatum, vmPFC, PCC and ACC. (C) Left panel: Regions encoding both the standard PE term and the difference term to biased PEs (conjunction) at different cluster-forming thresholds (color). Clusters significant at a threshold of z > 3.1 are surrounded by black edges. In bilateral striatum, vmPFC, bilateral ITG, and primary visual cortex, BOLD is significantly better explained by biased learning than by standard learning. Clusters in ACC, left motor cortex, and PCC are not significant any more.

1241

1242    Regarding EEG findings in this subgroup, both midfrontal theta and beta power reflected
1243    outcome valence: Theta power was higher for unfavorable than favorable outcomes (driven by a
1244    cluster around 225–500 ms, $p$ = .002), while beta power was higher for favorable than unfavorable
1245    outcomes (driven by a cluster around 325–1000 ms, $p$ = .002). When using PE terms as regressor for
1246    midfrontal EEG power while controlling for PE valence, delta power did not encode $PE_{STD}$ positively,
1247    though not significant (p = .056), and also the positive encoding of $PE_{DIF}$ was non-significant ($p$ =
1248    .053). The positive correlation of beta power with $PE_{STD}$ was not significant anymore ($p$ = .059),
1249    while the negative correlation with $PE_{DIF}$ remained (p = .001, 450–950 ms). When adding $PE_{STD}$ and
1250    $PE_{DIF}$ together to achieve $PE_{BIAS}$, theta/delta power indeed significantly encoded $PE_{BIAS}$, first

1251  positively ($p$ = .032, 224–475 ms) and then negatively ($p$ = .019, 600 – 1,000 ms; around 8 Hz and thus
1252  rather in the alpha band). Also, beta power was significantly negatively correlated with $PE_{BIAS}$ ($p$ =
1253  .008, 450 – 975 ms).
1254      In sum, all findings reported in the main text also held when analyzing only this subgroup of
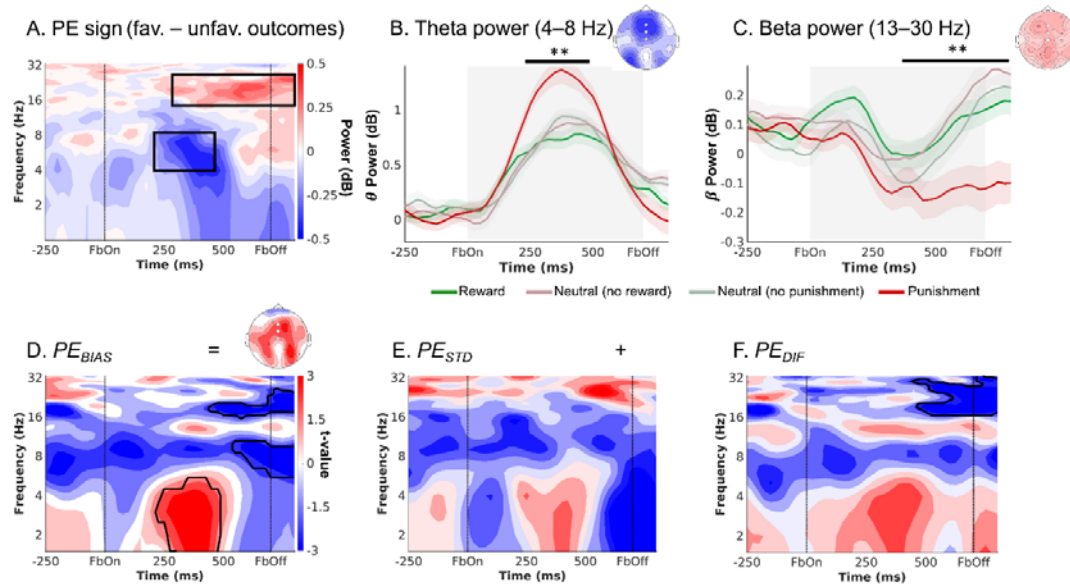1255  only 29 participants. In addition, also late beta power and theta/alpha power appeared to negatively
1256  encode the $PE_{BIAS}$ term.



**Figure S01C. EEG time-frequency power midfrontal electrodes (Fz/ FCz/ Cz) reflecting outcomes processing in the subgroup of 29 participants included in the fMRI-inspired EEG analyses.** (A) Time-frequency plot (logarithmic y-axis) displaying high theta (4–8 Hz) power for unfavorable outcomes and higher beta power (16–32 Hz) for favorable outcomes. (B). Theta power transiently increases for any outcome, but more so for unfavorable outcomes (especially punishments) around 225–475 ms after feedback onset. (C) Beta is higher for favorable than unfavorable outcomes (especially punishments) over a long time period around 300–1,250 ms after feedback onset. (D-F). Correlations between midfrontal EEG power and trial-by-trial PEs. Solid black lines indicate clusters above threshold. Biased PEs were significantly positively correlated with midfrontal theta power, but also negatively correlated with later alpha and beta power (D). The correlations of theta with the standard PEs (E) and the difference term to biased PEs (F) were also positive, though not significant. Beta power only encoded the difference term to biased PEs (F). ** p < 0.01. ** p < 0.01.

1257
1258      Regarding fMRI correlates of the past action, similar to the original analysis comprising 34
1259  participants, there were no clusters with higher BOLD after Go than NoGo actions at the time of
1260  outcomes, but vice versa, large parts of cortex and subcortex showed higher BOLD after NoGo than
1261  Go actions, highly similar to the original analysis ($z_{max}$ = 7.65, $p$ = 0, 124629 voxels, xyz = [-58 18 22]).
1262      Furthermore, there were four clusters with higher BOLD for Go than NoGo actions at the
1263  time of the response, namely one large cluster across lateral prefrontal cortex, anterior cingulate
1264  cortex, striatum, thalamus, angular gyrus, cerebellum, left operculum and motor cortex,
1265  intracalcarine cortex, and occipital pole ($z_{max}$ = 7.45, $p$ = 0, 61057 voxels, xyz = [32 -4 -4]), one in right
1266  middle temporal gyrus ($z_{max}$ = 4.90, $p$ = 8.66e-05, 493 voxels, xyz = [66 -32 -12]), one in left inferior
1267  temporal gyrus ($z_{max}$ = 4.43, $p$ = .00294, 293 voxels, xyz = [-60 -44 -18]), and one in precuneous ($z_{max}$ =
1268  2.39, $p$ = .0041, 276 voxels, xyz = [-8 -70 38]). All these regions were also found in the original analysis
1269  comprising 34 participants. Vice versa, BOLD signal was higher NoGo than Go actions at the time of
1270  the response in two clusters in vmPFC and subcallosal cortex ($z_{max}$ = 4.23, $p$ = .00864, 239 voxels, xyz

1271   = [-2 18 -6]) and right anterior temporal gyrus/ temporal pole ($z_{max}$ = .4.14, $p$ = .0193, 201 voxels, xyz
1272   = [48 -6 -8]), identical to the original analysis comprising 34 participants.
1273          Finally, there was higher BOLD signal for left hand compared to right hand responses at the
1274   time of response in two clusters in right precentral and postcentral gyrus, superior parietal lobule,
1275   and operculum ($z_{max}$ = 6.66, $p$ = 0, 11597 voxels, xyz = [46 -24 64]) and left cerebellum ($z_{max}$ = 6.76, $p$ =
1276   1.05e-18, 2672 voxels, xyz = [-18 -54 -16]), identical to the original analysis comprising 34
1277   participants. Vice versa, there was higher BOLD signal for right hand than left hand responses at the
1278   time of responses in five clusters in left precentral and postcentral gyrus, superior parietal lobule,
1279   operculum, and thalamus ($z_{max}$ = 6.4, $p$ = 0, 12372 voxels, xyz = [-36 -20 66]), right cerebellum ($z_{max}$ =
1280   7.17, $p$ = 3.41e-21, 3206 voxels, xyz = [20 -54 -20]), right superior lateral occipital cortex ($z_{max}$ = 4.84,
1281   $p$ = 2.28e-09, 988 voxels, xyz = [48 -86 -4]), right angular gyrus ($z_{max}$ = 4.11, $p$ = 7.68e-05, 396 voxels,
1282   xyz = [66 -50 28]), and left superior lateral occipital cortex ($z_{max}$ = 5.03, $p$ = .019, 164 voxels, xyz = [-18
1283   -82 48]). The clusters in right occipital pole/ intracalcarine cortex and in right posterior cerebellum
1284   observed in the original analysis comprising 34 participants were not observed in this analysis. In
1285   sum, all major findings also held when analyzing only this subgroup of only 29 participants.
1286          Regarding EEG time-frequency correlates of the past action, when testing for differences in
1287   broadband after outcome onset, there was no significant difference after Go and NoGo responses, $p$
1288   = .283. When restricting analyses to the low alpha range, the permutation test was marginally
1289   significant, $p$ = .056, driven by a cluster around 0–100 ms around 7–10 Hz). When repeating the
1290   permutation test for the broadband signal including the last second before outcome onset, there was
1291   a significant difference after Go and NoGo responses, driven by clusters in the beta band. $p$ = 0.002, -
1292   1000 − -275 ms, 13–32 Hz, and in the theta/ low alpha band, $p$ = 0.020, -1000 − -525 ms, 4–10 Hz.
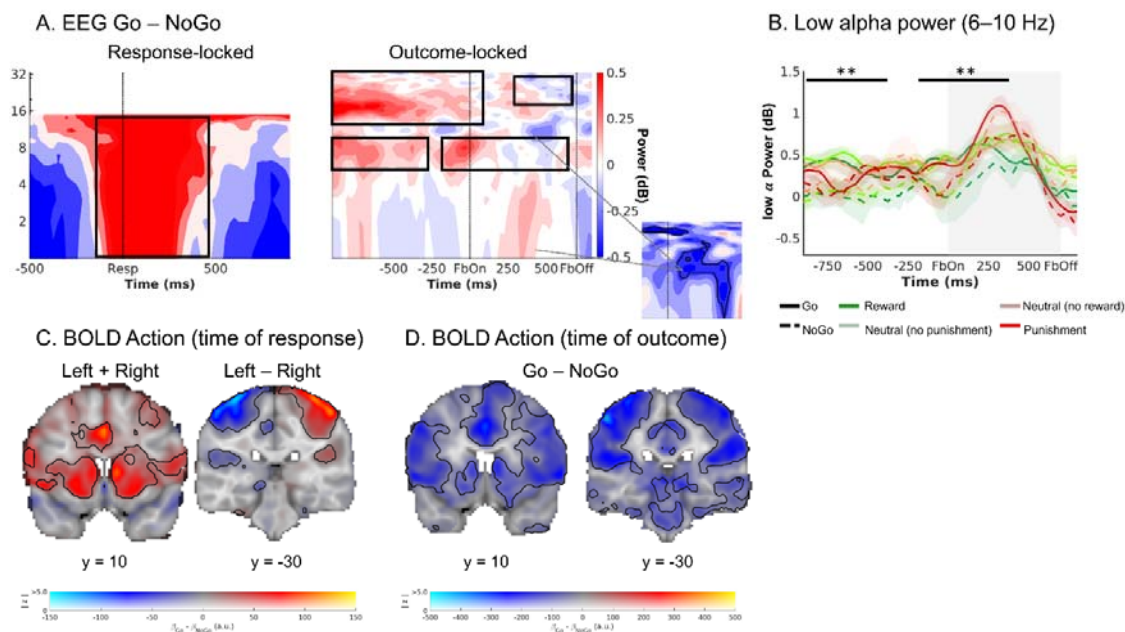


Figure S01D. Exploratory follow-up analyses on ACC BOLD signal and midfrontal low-alpha power in the subgroup of 29 participants included in the fMRI-inspired EEG analyses. (A) Midfrontal time-frequency response-locked (left panel) and outcome-locked (right panel). Before and shortly after outcome onset, power in the lower alpha band is higher on trials with Go actions than on trials with NoGo actions. The shape of this difference resembles the shape of ACC BOLD-EEG TF correlations (small plot; note that this plot depicts BOLD-EEG correlations, which are negative). Note that differences between Go and NoGo trials occurred already before outcome onset in the alpha and beta range, reminiscent of delay activity; but were not fully sustained since the actual response. (B) Midfrontal power in the lower alpha band per action x outcome condition. Lower alpha band power is consistently higher on trials with Go actions than on trials with NoGo

actions, starting already before outcome onset. (C) BOLD signal differences between Go and NoGo actions (left panel) and left vs. right hand responses (right panel) at the time or responses. Response-locked ACC BOLD is significantly higher for Go than NoGo actions. (D) BOLD signal differences between Go and NoGo actions at the time of outcomes. Outcome-locked ACC BOLD (and BOLD in other parts of cortex) is significantly lower on trials with Go than on trials with NoGo actions.

1293

1294    When linking trial-by-trial BOLD signal in selected ROIs as well as midfrontal EEG TF power to
1295 response repetition on the next trial with the same cue, ACC BOLD signal did not significantly predict
1296 the response repetition, $b$ = -0.013, $SE$ = 0.018, $\chi^2(1)$ = 0.524, $p$ = .469, and neither did PCC BOLD
1297 signal, $b$ = -0.037, $SE$ = 0.018, $\chi^2(1)$ = 2.079, $p$ = .149. However, participants in this subgroup were
1298 significantly more likely to repeat the sample action when striatal BOLD signal was high, $b$ = 0.097, $SE$
1299 = 0.025, $\chi^2(1)$ = 12.043, $p$ < .001, but more likely to switch when vmPFC BOLD was high, $b$ = -0.075, $SE$
1300 = 0.019, $\chi^2(1)$ = 13.170, $p$ < .001.
1301    When linking trial-by-trial midfrontal EEG TF power to response repetition on the next trial
1302 with the same cue, participants in this subgroup were more likely to repeat the same response when
1303 beta power was high, $b$ = 0.124, $SE$ = 0.036, $\chi^2(1)$ = 3.502, $p$ < .001, or when low alpha power was
1304 high, $b$ = 0.135, $SE$ = 0.044, $\chi^2(1)$ = 8.789, $p$ = .003, but more likely to switch to another response
1305 when theta power was high, $b$ = -0.090, $SE$ = 0.040, $\chi^2(1)$ = 4.812, $p$ = .028.

1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331

1332    **S02: Stay behavior as a function of action, salience, and valence**

1333

| Effect | $\chi^2$ | Df | *p*-value |
|---|---|---|---|
| Action | 0.01 | 1 | .924 |
| Salience | 5.15 | 1 | .021 |
| Valence | 45.59 | 1 | < .001 |
| Action x Salience | 0.12 | 1 | .728 |
| Action x Valence | 3.24 | 1 | .067 |
| Salience x Valence | 30.95 | 1 | < .001 |
| Action x Valence x Salience | 19.73 | 1 | < .001 |
| | | | |
| *Salient outcomes only:* | | | |
| Action | 0.01 | 1 | .960 |
| Valence | 46.36 | 1 | < .001 |
| Action x Valence | 17.80 | 1 | < .001 |
| | | | |
| *Neutral outcomes only:* | | | |
| Action | .102 | 1 | .750 |
| Valence | .830 | 1 | .362 |
| Action x Valence | 12.32 | 1 | < .001 |
| | | | |
| *Go with salient outcomes only:* | | | |
| Valence | 53.93 | 1 | < .001 |
| *NoGo with salient outcomes only:* | | | |
| Valence | 18.23 | 1 | < .001 |
| *Go with neutral outcomes only:* | | | |
| Valence | 0.13 | 1 | .050 |
| *NoGo with neutral outcomes only:* | | | |
| Valence | 7.21 | 1 | .007 |

**Table S02. Full report of model of stay behavior.** Mixed-effects logistic regression of stay vs. switch behavior (i.e., repeating vs. changing an action on the next occurrence of the same cue) as a function of performed action (Go vs. NoGo), outcome salience (salient: reward or punishment vs. neutral: no reward or no punishment), and outcome valence (positive: reward or no punishment vs. negative: no reward or punishment). Follow-up analyses are performed on trials with salient vs. neutral outcomes separately, and then separately based on Go vs. NoGo actions and salient vs. neutral outcomes. *P*-values are computed using likelihood ratio tests using the *mixed*-function (option "LRT") from package *afex*.

# S03: Model parameters and fit indices for models M1-M6

| | M1 | M2 | M3 | M4 | M5 (Asymmetric pathways) | M6 (Action priming) |
|---|---|---|---|---|---|---|
| Mean log model evidence | -609.30 | -597.95 | -554.46 | -532.40 | -528.13 | -540.84 |
| Model frequency | 0 | 0.0278 | 0 | 0.0488 | 0.6815 | 0.2419 |
| Protected exceedance probability | 0 | 0 | 0 | 0 | .9970 | .0030 |
| $\rho$ | 7.75 $[0.53 - 38.68]$ | 6.81 $[0.48 - 37.74]$ | 6.38 $[0.49 - 35.71]$ | 10.05 $[1.26 - 40.60]$ | 9.41 $[0.98 - 31.22]$ | 6.64 $[0.71 - 22.83]$ |
| $\varepsilon_0$ | 0.17 $[0.002 - 0.77]$ | 0.20 $[0.003 - 0.82]$ | 0.21 $[0.003 - 0.85]$ | 0.09 $[0.003 - 0.38]$ | 0.08 $[0.003 - 0.41]$ | 0.039 $[0.003 - 0.11]$ |
| b | | -0.05 $[-1.23 - 0.82]$ | -0.01 $[-1.23 - 1.09]$ | 0.13 $[-1.16 - 1.03]$ | 0.14 $[-1.18 - 1.10]$ | 0.16 $[-1.22 - 1.40]$ |
| $\pi$ | | | 0.77 $[-0.78 - 3.73]$ | | 0.17 $[-1.25 - 2.70]$ | -1.11 $[-3.29 - 1.23]$ |
| $\varepsilon_{\text{rewarded Go}}\,(\varepsilon_0+\kappa)$ | | | | 0.749 $[0.29 - 0.99]$ | 0.833 $[0.43 - 0.99]$ | |
| $\varepsilon_{\text{punished NoGo}}\,(\varepsilon_0-\kappa)$ | | | | 0.001 $[0.001 - 0.02]$ | 0.003 $[0.001 - 0.09]$ | |
| $\varepsilon_{\text{salient Go}}$ | | | | | | 0.49 $[0.05 - 0.90]$ |

**Table S03. Model parameters for fitted models.** Mean [minimum − maximum] of participant-level parameter estimates in model space, fitted with hierarchical Bayesian inference (only the respective model included in the fitting process). Model frequency and protected exceedance probability are based on a model comparison that involves models M1-M6. Note that Fig. 2 in the main text does not include M6.

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

## S04: Simulations for asymmetric pathways and action priming model

Motivational learning biases are predicted by the *asymmetric pathways model* (*15*, *16*): Positive PEs, elicited by rewards, lead to long-term potentiation in the striatal direct "Go" pathway (and long term depression in the indirect pathway), allowing for a particularly effective acquisition of Go actions to obtain rewards. Conversely, negative PEs, elicited by punishments, lead to long term potentiation in the NoGo pathway, impairing the unlearning of NoGo actions in face of punishments.

An alternative account has recently suggested that self-generated (Go) actions lead to preferential learning (relative to non-self-generated actions, including inaction), more generally (henceforth called "action priming model")(*33*). A self-generated action could "prime" basal ganglia circuits and lead to subsequently larger PEs and thus faster learning. The main differential prediction between these two models is how they account for the failure to learn "Go" actions to avoid punishment: In the first model, this is due to a failure to unlearn punished "NoGo" actions, while in the second model, this is due increased unlearning of punished "Go" actions.

Here, we directly tested both models against each other. As an alternative model M6 (Cockburn et al. 2014), we specified a model with two separate learning rates, one learning rate for trials where self-generated (Go) action selection should prime the processing of any following salient outcome (i.e., Go actions followed by rewards/ punishments), and one learning rate for any other action-outcome combination. In this model, equation (6) is substituted by equation (7):

$$\varepsilon = \begin{cases} \varepsilon_{salGo} \; for \; any \; Go \; action \; with \; salient \; outcomes \\ \varepsilon_0 \hspace{4.5cm} else \end{cases} \quad (7)$$

When comparing all models M1–M6 using Bayesian model selection, M5 (the asymmetric pathways model) received highest support (model frequency: 68.15%; protected exceedance probability: 99.70%), also compared to M6 (the action priming model; model frequency: 24.19%; protected exceedance probability: 0.30%). In fact, as visible in Fig. S04E-H, the action priming did not reproduce the motivational biases in learning curves and bar plots, which constitutes a case of qualitative model falsification (*79*, *80*). If anything, it seems that the action priming model trades off both biases, leading to negative response biases for a majority of participants. In contrast, the asymmetric pathways model (M5) was well able to capture the qualitative patterns observed in the data (Fig. S04A-D). We conclude that only the asymmetric pathways model is able to qualitatively reproduce core characteristics of our data.
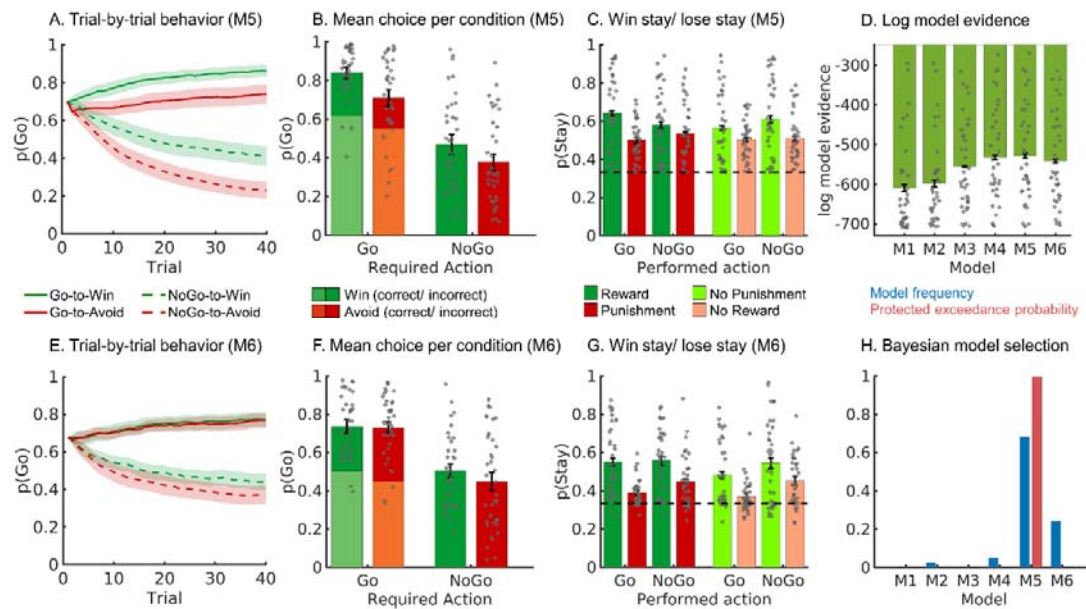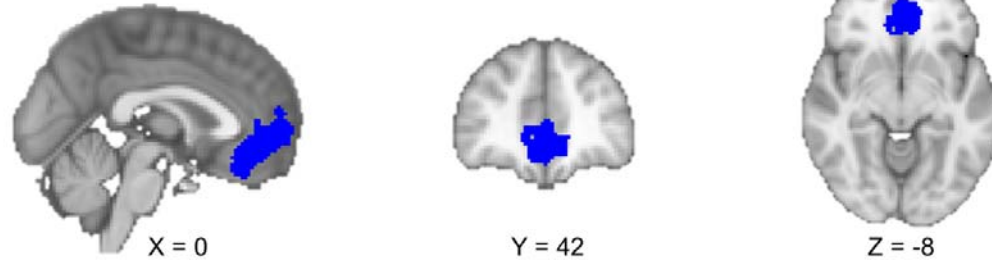
**Figure S04. Model comparison and validation of asymmetric pathways (M5) and action priming (M6) model.** (A-C) One-step-ahead predictions using parameters (hierarchical Bayesian inference) of the winning model asymmetric pathways model (M5). (A) Trial-by-trial proportion of Go responses (±SEM across participants) for Go cues (solid lines) and NoGo cues (dashed lines); (B) Mean (±SEM across participants) proportion Go responses per cue condition (points are individual participants' means); (C) Probability to repeat a response ("stay") on the next encounter of the same cue as a function of action and outcome. The asymmetric pathways model is well able to capture core characteristics of the empirical data (see Fig. 2 in the main text). (D) Log-model evidence favors the asymmetric pathways model (M5), even over the action priming model (M6). (E-G) Trial-by-trial proportion of Go responses, mean proportion Go responses, and probability of for the action priming model (M6). This model does not reproduce motivational biases (i.e., the difference between green and red lines and bars) well. (H) Model frequency and protected exceedance probability indicate best fit for model M5 (asymmetric pathways model), in line with log model evidence.
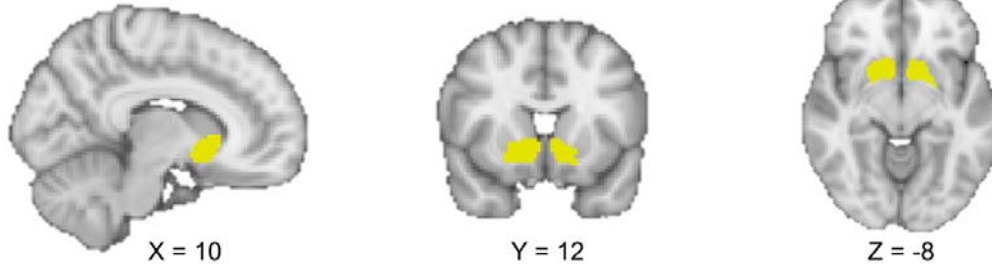
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426

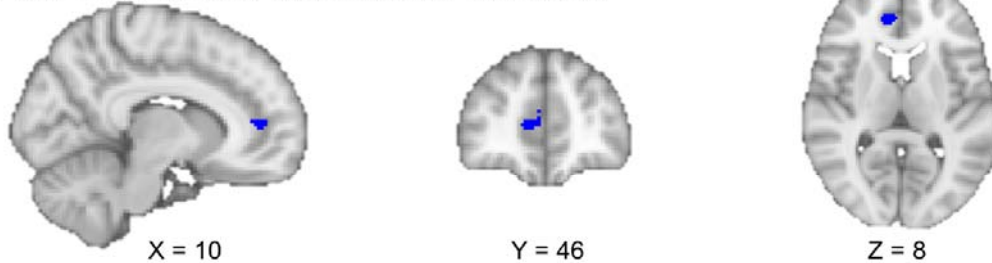1427  ## S05: Anatomical masks and conjunctions of anatomical and
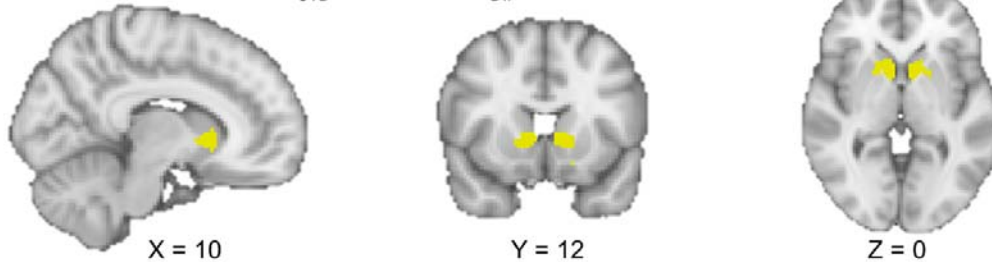
1428  ## functional masks

1429



**Figure S05A. Conjunctions of anatomical masks with functional contrasts from fMRI GLM analyses used for fMRI-informed EEG analyses.** Anatomical masks are based on the Harvard-Oxford Atlas. Functional contrasts involve outcome Valence and conjunction of $PE_{STD}$ and $PE_{DIF}$. (A) Anatomical AAC contrast (pink, cingulate gyrus, anterior division); (B) vmPFC outcome valence contrast (dark blue, conjunction of frontal pole, frontal medial cortex, and paracingulate gyrus); (C) striatum outcome valence contrast (yellow, conjunction of bilateral nucleus accumbens, caudate, and putamen); (D) vmPFC $PE_{STD} \cap PE_{DIF}$ contrast (dark blue); and (E) and striatum $PE_{STD} \cap PE_{DIF}$ contrast (yellow). All anatomical masks were extracted from the probabilistic Harvard-Oxford Atlas, thresholded at 10%. Note that images are in radiological orientation (i.e., left brain hemisphere presented on the right and vice versa).

1430



**Figure S05B. Conjunctions of anatomical masks with functional contrasts from fMRI GLM analyses used for fMRI-informed EEG analyses**: (A) AAC $PE_{STD} \cap PE_{DIF}$ contrast (red, cingulate gyrus, anterior division); (B) PCC $PE_{STD} \cap PE_{DIF}$ contrast (light blue, cingulate gyrus, posterior division); (C) left motor cortex $PE_{STD} \cap PE_{DIF}$ contrast (orange, conjunction of precentral and postcentral gyrus); (D) Left inferior temporal gyrus $PE_{STD} \cap PE_{DIF}$ contrast (turquoise, conjunction of inferior temporal gyrus, posterior division, and inferior temporal gyrus, temporooccipital part); and (E) primary visual cortex $PE_{STD} \cap PE_{DIF}$ contrast (green, conjunction of lingual gyrus, occipital fusiform gyrus, occipital pole). All anatomical masks were extracted from the probabilistic Harvard-Oxford Atlas, thresholded at 10%. Note that images are in radiological orientation (i.e., left brain hemisphere presented on the right and vice versa).
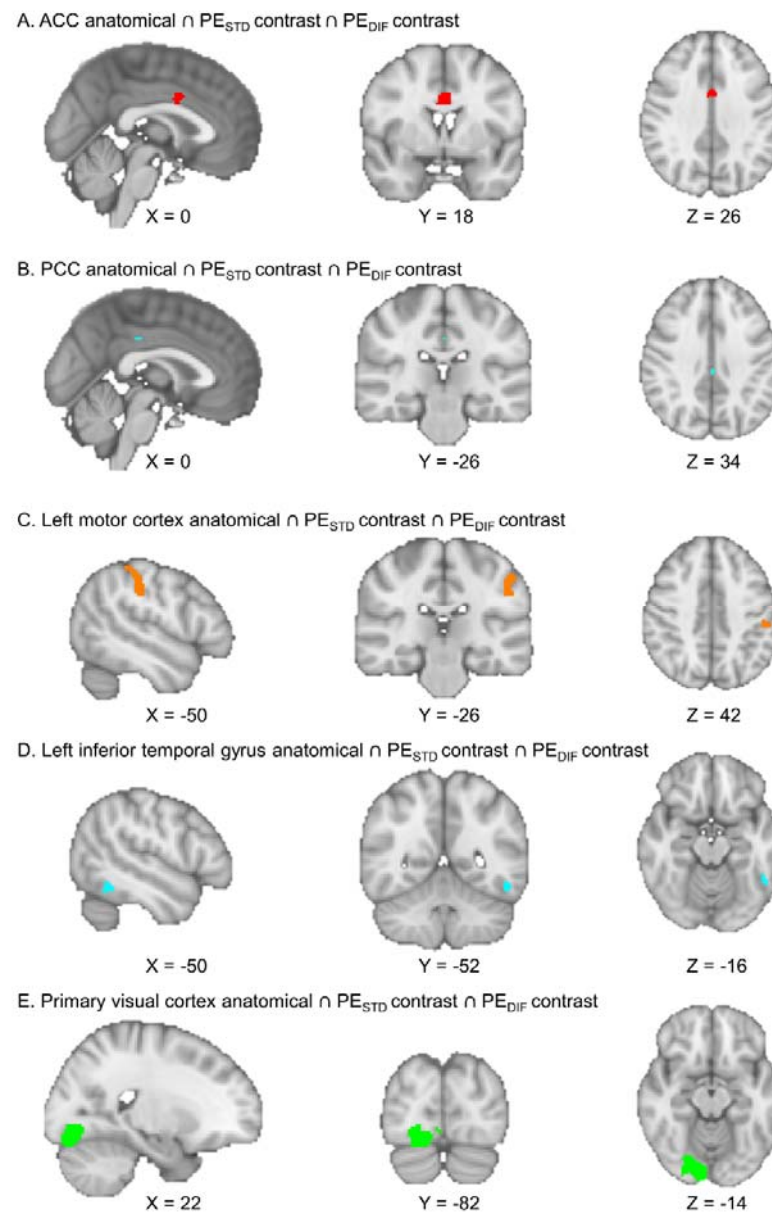
1431
1432
1433
1434

## S06: Regressors and contrast in fMRI analyses

**Model-based GLM with $PE_{STD}$ and $PE_{DIF}$ regressor:**

- WinGoOnset: for every trial with Win cue and Go action, at cue onset, duration 1, value +1
- AvoidGoOnset: for every trial with Avoid cue and Go action, at cue onset, duration 1, value +1
- WinNoGoOnset: for every trial with Win cue and NoGo action, at cue onset, duration 1, value +1
- AvoidNoGoOnset: for every trial with Avoid cue and NoGo action, at cue onset, duration 1, value +1
- Handedness: for every trial, at cue onset, duration 1, value +1 for left hand response, 0 for NoGo 10 response, -1 for right hand response 11
- Error: for every trial, at cue onset, duration 1, value +1 for incorrect response, 0 for correct response
- OutcomeOnset: for every trial, at outcome onset, duration 1, value +1 for every trial
- $PE_{STD}$: for every trial, at outcome onset, duration 1, value is demeaned PE times learning rate for model M1
- $PE_{DIF}$: for every trial, at outcome onset, duration 1, value is demeaned difference between (PE times learning rate) for model M1 and (PE times learning rate) for model M5
- Invalid: for trials where uninstructed button was pressed, at outcome onset, duration 1, value 1

| | Contrast | 1 WinGoOnset | 2 AvoidGoOnset | 3 WinNoGoOnset | 4 AvoidNoGoOnset | 5 Handedness | 6 Error | 7 Outcome Onset | 8 $PE_{STD}$ | 9 $PE_{DIF}$ | 10 Invalid |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | $PE_{STD}$ | | | | | | | | 1 | | |
| 2 | $PE_{DIF}$ | | | | | | | | | 1 | |

1471 **Model-free GLM using response-locked and outcome-locked response regressors:**

1472 • GoReward: for every trial with Go action and reward obtained, at outcome onset, duration 1,
1473 value +1

1474 • GoNoReward: for every trial with Go action and no reward obtained, at outcome onset,
1475 duration 1, value +1

1476 • GoNoPunishment: for every trial with Go action and no punishment obtained, at outcome
1477 onset, duration 1, value +1

1478 • GoPunishment: for every trial with Go action and punishment obtained, at outcome onset,
1479 duration 1, value +1

1480 • NoGoReward: for every trial with NoGo action and reward obtained, at outcome onset,
1481 duration 1, value +1

1482 • NoGoNoReward: for every trial with NoGo action and no reward obtained, at outcome onset,
1483 duration 1, value +1

1484 • NoGoNoPunishment: for every trial with NoGo action and no punishment obtained, at
1485 outcome onset, duration 1, value +1

1486 • NoGoPunishment: for every trial with NoGo action and punishment obtained, at outcome
1487 onset, duration 1, value +1

1488 • LeftHand: for very trial with left hand response, at response onset, duration 1, value + 1

1489 • RightHand: for very trial with right hand response, at response onset, duration 1, value + 1

1490 • Error: for every trial, at cue onset, duration 1, value +1 for incorrect response, 0 for correct
1491 response

1492 • OutcomeOnset: for every trial, at outcome onset, duration 1, value +1 for every trial

1493 • Invalid: for trials where uninstructed button was pressed, at outcome onset, duration 1,
1494 value 1

1495

| Regressors | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Contrast | GoReward | GoNoReward | GoNoPunishment | GoPunishment | NoGoReward | NoGoNoReward | NoGoNoPunishment | NoGoPunishment | LeftHand | RightHand | Error | OutcomeOnset | Invalid |
| 1 | Valence | 1 | -1 | 1 | -1 | 1 | -1 | 1 | -1 | | | | | |
| 2 | Action | 1 | 1 | 1 | 1 | -1 | -1 | -1 | -1 | | | | | |
| 3 | Hand Sum | | | | | | | | | 1 | 1 | | | |
| 4 | Hand Dif | | | | | | | | | 1 | -1 | | | |

1496
1497
1498
1499
1500
1501
1502

## S07: Significant clusters in BOLD-GLMs with behavioral regressors only

**Model-based GLM with PE$_{STD}$ and PE$_{DIF}$ regressor:**

| No | Contrast / Brain region | Maximal Z-value | Cluster size (voxels) | Corrected p | Peak coordinates x | y | z |
|---|---|---|---|---|---|---|---|
| | **PE$_{STD}$ Positive** | | | | | | |
| 1 | Ventromedial prefrontal cortex, Nucleus accumbens, caudate, putamen, bilateral amygdala, bilateral hippocampus | 6.47 | 8762 | 1.02e-43 | 12 | 14 | -6 |
| 2 | Occipital pole, lingual gyrus, occipital fusiform gyrus | 6.64 | 1012 | 6.10e-10 | 10 | -92 | -10 |
| 3 | Posterior cingulate cortex | 4.72 | 985 | 9.40e-10 | 4 | -50 | 18 |
| 4 | Left superior frontal gyrus | 5.56 | 910 | 3.19e-09 | -18 | 34 | 50 |
| 5 | Right middle temporal gyrus, anterior division | 5.48 | 381 | 6.47e-05 | 62 | -4 | -18 |
| 6 | Left inferior temporal gyrus, temporooccipital part | 5.16 | 360 | .000103 | -52 | -46 | -10 |
| 7 | Left middle temporal gyrus, anterior division | 4.70 | 329 | .000209 | -60 | -10 | -14 |
| 8 | Left postcentral gyrus | 4.33 | 271 | .000838 | -52 | -28 | 48 |
| 9 | Right cerebellum | 4.89 | 147 | .0239 | 44 | -72 | -40 |
| 10 | Anterior cingulate cortex | 4.27 | 146 | .0247 | 2 | 6 | 34 |
| | **PE$_{STD}$ Negative** | | | | | | |
| 1 | Right superior frontal gyrus | 5.20 | 351 | .000127 | 6 | 26 | 62 |
| 2 | Right occipital pole, right inferior lateral occipital cortex | 4.76 | 211 | .00391 | 30 | -94 | 4 |
| 3 | Left lingual gyrus | 4.21 | 186 | .00776 | -22 | -64 | 2 |
| 4 | Left inferior lateral occipital cortex | 4.28 | 147 | .0239 | -44 | -86 | -10 |
| | **PE$_{DIF}$ Positive** | | | | | | |
| 1 | Bilateral superior frontal gyrus, paracingulate gyrus, anterior cingulate cortex, posterior cingulate cortex, ventromedial frontal cortex, bilateral frontal orbital cortex, bilateral frontal pole, bilateral supramarginal gyrus, bilateral middle temporal gyrus, bilateral inferior temporal gyrus, bilateral fusiform gyrus, bilateral inferior occipital cortex, bilateral superior occipital cortex, precuneous, bilateral cerebellum | 7.11 | 35109 | 0 | 34 | -84 | 20 |
| 2 | Right insula, right frontal operculum, right inferior frontal gyrus, right middle frontal gyrus, right frontal orbital cortex, bilateral caudate, bilateral Nucleus accumbens, bilateral thalamus, brainstem | 6.36 | 10364 | 0 | 34 | 20 | -8 |
| 3 | Left insula, left frontal operculum, | 6.51 | 10132 | 0 | -36 | 20 | -6 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | left inferior frontal gyrus, left middle frontal gyrus, left frontal orbital cortex | | | | | | |
| 4 | Right middle temporal gyrus, posterior division | 4.66 | 307 | .0003 | 56 | -32 | -4 |
| 5 | Right insula, right planum polare | 4.72 | 143 | .0248 | 40 | -8 | -12 |
| | **PE$_{DIF}$ Negative** | | | | | | |
| 1 | Left middle temporal gyrus, anterior division | 4.22 | 191 | .00607 | -64 | -6 | -14 |
| 2 | Left hippocampus | 4.49 | 158 | .0158 | -26 | -14 | -22 |

1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543

1544     **Model-free GLM using response-locked and outcome-locked response regressors:**

1545

| | Contrast | | | | Peak coordinates | | |
|---|---|---|---|---|---|---|---|
| No | Brain region | Maximal Z-value | Cluster size (voxels) | Corrected p | x | y | z |
| | **Favorable > Unfavorable** | | | | | | |
| 1 | Ventromedial prefrontal cortex, left lateral orbitofrontal cortex, Nucleus accumbens, caudate, putamen, bilateral amygdala, bilateral hippocampus | 5.65 | 3999 | 2.86e-19 | 8 | 12 | -4 |
| 2 | Left superior frontal gyrus | 4.03 | 331 | 0.00239 | -18 | 28 | 60 |
| 3 | Left lateral orbitofrontal cortex | 4.31 | 288 | 0.00512 | -34 | 40 | -8 |
| 4 | Right occipital pole | 4.59 | 213 | 0.0212 | 18 | -92 | -16 |
| | **Unfavorable > Favorable** | | | | | | |
| 1 | Right lateral orbitofrontal cortex | 4.59 | 367 | 0.00142 | 30 | 62 | -2 |
| 2 | Precuneous | 4.58 | 356 | 0.00170 | 8 | -66 | 58 |
| 3 | Right superior frontal gyrus | 4.32 | 340 | 0.00223 | 12 | 14 | 72 |
| | **Go > NoGo** **outcome-locked** *No significant clusters* | | | | | | |
| | **NoGo > Go** **outcome-locked** | | | | | | |
| 1 | Bilateral lateral orbitofrontal cortex, Bilateral superior frontal gyrus, anterior cingulate cortex, posterior cingulate cortex, pre-SMA, bilateral precentral gyrus, bilateral postcentral gyus, bilateral supramarginal gyrus, bilateral operculum, bilateral planum temporale, bilateral superior temporal gyrus, bilateral middle temporal gyrus, bilateral inferior temporal gyrus, bilateral superior lateral occipital cortex, bilateral inferior lateral occipital cortex, bilateral thalamus | 7.32 | 114090 | 0 | -42 | -6 | 12 |
| | **Go (left + right hand response) > NoGo** **response-locked** | | | | | | |
| 1 | Cerebellum, bilateral thalamus, bilateral putamen, bilateral caudate, bilateral Nucleus Accumbens, posterior cingulate cortex, right operculum, right angular gyrus, right superior parietal lobule. anterior cingulate cortex, paracingulate gyrus, bilateral ventrolateral frontal cortex, right middle frontal gyrus | 7.08 | 46437 | 0 | 32 | -4 | -6 |
| 2 | Left operculum, left angular gyrus, left superior parietal lobule | 5.88 | 3936 | 3.13e-17 | -46 | -24 | 26 |
| 3 | Intracalcarine cortex | 3.79 | 374 | 0.00248 | -12 | -88 | 6 |
| 4 | Right middle temporal gyrus | 4.63 | 287 | 0.00956 | 68 | -32 | -12 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | **NoGo > Go (left + right hand response) response-locked** | | | | | | |
| 1 | Right medial temporal gyrus, right temporal pole | 4.09 | 465 | 0.000636 | 50 | -8 | -16 |
| 2 | vmPFC, subcallosal cortex | 3.95 | 435 | 0.000973 | 0 | 40 | -12 |
| | **Left Hand > Right Hand Response response-locked** | | | | | | |
| 1 | Right precentral gyrus, right postcentral gyrus, right superior parietal lobule, right operculum | 7.05 | 9460 | 9.41e-39 | 46 | -24 | 64 |
| 2 | Left cerebellum | 7.18 | 2208 | 2.1e-14 | -18 | -54 | -18 |
| | **Right Hand > Left Hand Response response-locked** | | | | | | |
| 1 | left precentral gyrus, left postcentral gyrus, left superior parietal lobule, left operculum, left thalamus | 7.06 | 14870 | 0 | -36 | -20 | 66 |
| 2 | Right anterior cerebellum | 7.90 | 3735 | 1.44e-20 | 18 | -54 | -20 |
| 3 | Right inferior lateral occipital cortex, right superior lateral occipital cortex | 4.96 | 1452 | 9.66e-11 | 48 | -86 | -4 |
| 4 | Right angular gyrus | 4.98 | 551 | 2.06e05 | 66 | -50 | 28 |
| 5 | Left occipital pole, right intracalcarine cortex | 3.93 | 409 | 0.000236 | -4 | -96 | 26 |
| 6 | Right posterior cerebellum | 4.64 | 200 | 0.0157 | 48 | -78 | -32 |

1546

1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1558

1559

1560

1561

1562

1563

1564

1565

1566

1567

1568

## S08: EEG time-frequency results after ERPs are removed

1569

1570 Given that differences in theta power between favorable and unfavorable outcomes as well
1571 as differences in lower alpha band power after Go and NoGo responses occurred quite soon after cue
1572 onset, we aimed to test whether these effects reflected differences in evoked rather than induced
1573 activity. For this purpose, we removed evoked components from our data by computing the ERP for
1574 each of the eight conditions (action x outcome) for each participant and then subtracting the
1575 condition-specific ERP from the trial-by-trial data (*81*). Only afterwards, we performed time-
1576 frequency decomposition.

1577 In line with the results reported in the main text, power was higher for unfavorable
1578 compared to favorable outcomes in the theta band ($p$ =.018, driven by cluster at 225–475 ms; Fig.
1579 S08B), but higher for favorable than unfavorable outcomes in the beta band ($p$ < .001, driven by
1580 cluster at 0–1250 ms; Fig. S08C). Notably, unlike the results reported in the main text (Fig. 4A), the
1581 cluster of high power for unfavorable compared to favorable outcomes was constrained to the theta
1582 range, and did not extend further into the delta range (Fig. S08A).

1583 When using the trial-by-trial PEs (both the standard PE and the difference term to a biased
1584 PE) as predictors in a multiple linear regression at each time-frequency-channel bin while controlling
1585 for PE valence, delta power encoded $PE_{STD}$ positively, though not significantly ($p$ = .198). However,
1586 at a later time point around outcome offset, delta (and theta) power in fact correlated negatively
1587 with $PE_{STD}$ (575–800 ms, $p$ = .002; Fig. S08E). The correlation between delta and the $PE_{DIF}$ term was
1588 still positive, but not significant (p = .228, Fig. S08F). Similarly, the correlation of the $PE_{BIAS}$ term
1589 with delta power was positive, but not significant ($p$ = .084; Fig. S08D).

1590 Regarding beta power, there was a positive, though non-significant correlation of beta power
1591 with $PE_{STD}$ ($p$ = .096). There was again a significantly negative correlation of beta power with $PE_{DIF}$
1592 (425–875 ms, $p$ < .001, Fig. S08B). Likewise, beta power correlated significantly negatively with
1593 $PE_{BIAS}$ (450–800 ms, $p$ = .018), driven by the correlation with $PE_{DIF}$.

1594 In sum, after subtracting the condition-wise ERP from each trial before time-frequency
1595 decomposition, supposedly removing the phase-locked aspect of power, both beta and theta still
1596 encoded PE valence. However, the encoding of PE magnitude by delta power was attenuated and not
1597 significant any more.

1598 This reduction in magnitude encoding might occur of several reasons. Firstly, it might be that
1599 this correlation in the delta range is in fact (partly) reflecting correlations with phase-locked, i.e.,
1600 evoked activity (ERPs), especially in the N2 (FPN)/ P3 (RewP) time range (see S09) (*26, 28–30, 82–87*).
1601 Nonetheless, a positively correlation between delta power and biased PEs is still visible in Fig. S08D,
1602 suggesting that at least part of the signal encoding biased PEs is not phase-locked. Secondly, it might
1603 be that the removal of the condition-wise ERPs has introduced additional noise in the data,
1604 attenuating any true correlation. Thirdly, there was a negative correlation between $PE_{STD}$ and theta/
1605 delta power at later time points which was visible, though not significant in the results reported in
1606 the main text (Fig. 4D). Subtraction of an ERP-like template acts like a high-pass filter. High-pass
1607 filtering at relatively high cut-offs (> 0.5 Hz) can artificially postpone or induce effects at later points
1608 (*88*). It is possible that in this case, ERP subtraction attenuated a positive correlation in the theta/
1609 delta range, but enhanced a later negative correlation.

1610 Taken together, it is possible that part of the PE magnitude encoding in the theta/ delta
1611 range is due to correlations with the phase-locked (ERP) signal. However, this finding does not
1612 compromise the conclusion that overall, theta/delta power seemed to be more strongly associated
1613 with the $PE_{BIAS}$ term than the $PE_{STD}$ term. Our primary goal is not to pinpoint the precise nature of

1614  electrophysiological correlates of biased learning, but rather test the relative temporal order of when
1615  different regions exhibiting biased learning signals become active.
1616          Finally, we tested whether after ERP subtraction, low alpha (and beta power) still encoded
1617  the previously performed action. When testing for differences in broadband power after Go and
1618  NoGo responses, power was indeed significantly different between conditions, driven by clusters in
1619  beta band ($p$ = 0.002, 0.125 − 625 ms; $p$ = 0.052, 700 - 1000 ms, 23 - 29 Hz) and theta/ low alpha
1620  band ($p$ = 0.024, 575 − 1000 ms, 5–9 Hz; $p$ = 0.056, 0 −225 ms, 6–11 Hz). For power before outcome
1621  onset, there were again broadband differences between Go and NoGo ($p$ = 0.002, -1000 − +225 ms, 1
1622  −33 Hz), but note that there was no ERP subtracted before outcome onset. We thus conclude that
1623  the differences between Go and NoGo responses were attributable to differences in induced rather
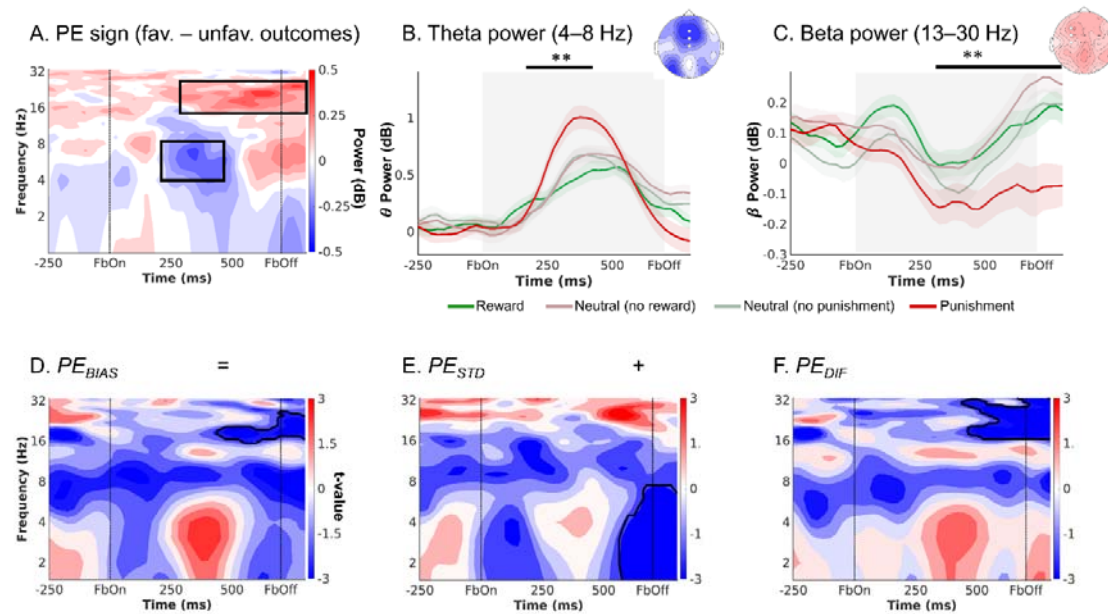1624  than evoked activity.
1625



**Figure S08. EEG time-frequency power over midfrontal electrodes (Fz/ FCz/ Cz) after the (action x outcome) condition-wise ERPs has been removed.** (A) Time-frequency plot (logarithmic y-axis) displaying high theta (4–8 Hz) power for unfavorable outcomes and higher beta power (16–32 Hz) for favorable outcomes. (B). Theta power transiently increases for any outcome, but more so for unfavorable outcomes (especially punishments) around 225–475 ms after feedback onset. (C) Beta is higher for favorable than unfavorable outcomes (especially punishments) over a long time period around 300–1,250 ms after feedback onset. (D-F). Correlations between midfrontal EEG power and trial-by-trial PEs. Solid black lines indicate clusters above threshold. There still was a visible positive correlation between biased PEs and midfrontal delta power, but this correlation was not significant (D). The correlation of delta with the standard PEs (E) was also positive, though not significant; in fact, at a later time point around stimulus offset, delta power correlated significantly negatively with standard PEs. The difference term to biased PEs (F) also correlated positively, though not significantly with delta power. Beta power encoded the difference term and biased PEs themselves (F). ** $p$ < 0.01.

1626
1627
1628
1629
1630

## S09: ERPs as a function of action and outcome

In addition to the induced activity in time-frequency power reported in the main text, we also analyzed the data in the time domain to test for differences in evoked activity. These analyses were particularly motivated given that differences in time-frequency power between favorable and unfavorable outcomes (theta/delta range) and after Go and NoGo responses (lower alpha/ theta range) occurred soon after outcome onset, warranting the assumption that differences might also occur in evoked activity. A large range of previous research has reported a modulation of evoked potentials by outcome valence in form of the feedback-reduced negativity (*29, 64, 82–87*), i.e., a stronger N2 component for negative compared to positive outcomes around ~ 250 post-cue over midfrontal electrodes, recently also characterized as rather constituting a reward positivity (RewP) (*82*). Also, some studies have reported a modulation of the P3 by outcome valence, which has been attributed to outcome magnitude or salience rather than valence (*85–87, 89*).

Similar to the analysis of time frequency power, we sorted trials into the eight conditions spanned by the performed action (Go/ NoGo) and the obtained outcome (reward/ no reward/ no punishment/ punishment), computed the average ERP for each condition per participant, and tested for differences between favorable (reward/ no punishment) and unfavorable (no reward/ punishment) outcomes as well as conditions of relative stronger (rewarded Go and punished Go) vs. relatively weaker learning (rewarded NoGo and punished NoGo). We used cluster-based permutation tests on the average signal over midfrontal electrodes (Fz/ FCz/ Cz) in the time range of 0–700 ms after outcome onset (where evoked potentials visible in condition-averaged plot).

First, midfrontal ERPs were significantly different between favorable and unfavorable outcomes, driven by two separate clusters of differences above threshold (Cluster 1: around 246 – 294 ms, $p$ = .034; Cluster 2: around 344 – 414 ms, $p$ =.004, Fig. S09A panel A, C). The first cluster the classical feedback-related negativity, i.e., a stronger N2 component for unfavorable compared to favorable outcomes. The second cluster reflected weaker P3 component for unfavorable compared to favorable outcomes, similar the reward positivity reported before. In fact, the N3 was rather absent for unfavorable outcomes (Fig. S09B). Both effects were clearly focused on midfrontal electrodes. These findings replicate previous findings of outcome valence modulating N2 (feedback-related negativity) and P3 components, and complement our time-frequency findings of theta and beta power reflecting outcome valence.

Second, when contrasting trials with Go vs. NoGo responses, no significant difference was observed ($p$ = .358; Fig. S09A panel D). Visual inspection of the topoplot yielded that, if anything, differences emerged over right occipital electrodes. If one performed a test over those right occipital electrodes (O2, 04, PO4; Fig. S09A panel F; note that this procedure constitutes double-dipping because the test was informed by first looking at the data), this test would have yielded significant results ($p$ = .016) driven by cluster around 423–466 ms, reflecting a slightly larger P3 after Go than NoGo responses (Fig. S09A panel E). This finding appears to be the strongest (if any) difference in amplitude after outcome onset between Go and NoGo actions. Given that this difference was not hypothesized and occurred far away from our a-priori selected channels of interest, we are careful not to over-interpret those differences.

Third, contrasting trials with favorable and unfavorable at the same right occipital electrodes yielded a significant difference, driven by clusters around 46–103 ms ($p$ = 0.034), 141–255 ms ($p$ = .002), and 519 – 580 ms ($p$ = .034). Most notably, the P1 amplitude was much larger for favorable than unfavorable outcomes (Fig. S09A panel B). However, given that these differences were not

1675     hypothesized and occurred far away from our a-priori selected channels of interest, we are careful
1676     not to over-interpret those differences.
1677         Taken together, we found a bigger midfrontal N2/ FRN for unfavorable compared to
1678     favorable outcomes, and a bigger midfrontal P3/ RewP for favorable compared to unfavorable
1679     outcomes, in line with a vast literature of previous findings (*29, 64, 82–87, 89*). Midfrontal voltage did
1680     not significantly differ after Go or NoGo responses. If anything, differences after Go and NoGo
1681     responses were maximal over right occipital electrodes, with a larger P3 after Go than after NoGo
1682     responses. Signal at these channels also differed between favorable and unfavorable outcomes, most
1683     notably with a bigger P1 after favorable than unfavorable outcomes. In sum, we replicate classical
1684     reward learning ERP effects, which shows that the motivational Go/NoGo learning task taps into
1685     reward learning processes reported before, but these processes appeared to be unaffected by the
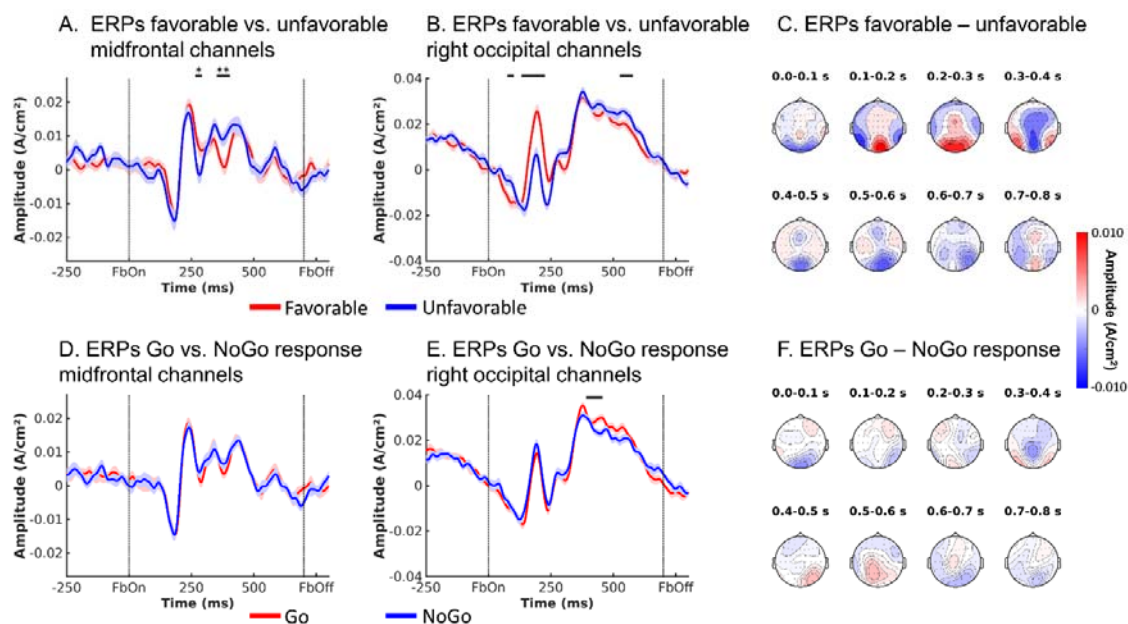1686     previously performed action.
1687



**Figure S09A. ERPs reflecting outcome valence and performed action**. (A) Voltage (±SEM) over midfrontal electrodes (Fz/FCz/Cz) was lower for unfavorable than favorable outcomes around 246–294 ms (stronger N2, FRN) and higher for favorable than unfavorable outcomes around 344 – 414 ms (stronger P3/ RewP). (B) Over right occipital electrodes, the P3 was slightly bigger for favorable than unfavorable outcomes. ** $p < 0.01$. * $p < .05$ (C) Topoplots of difference in voltage between trials with favorable and unfavorable outcomes over selected time windows. (D) There was no difference in voltage over midfrontal electrodes between trials with Go and NoGo responses. (E) Over right occipital electrodes, the P3 was slightly stronger after Go than NoGo actions (no *p*-value because ROI selected based on visual inspection). (F) Topoplots of difference in voltage between trials with Go and NoGo actions over selected time windows.
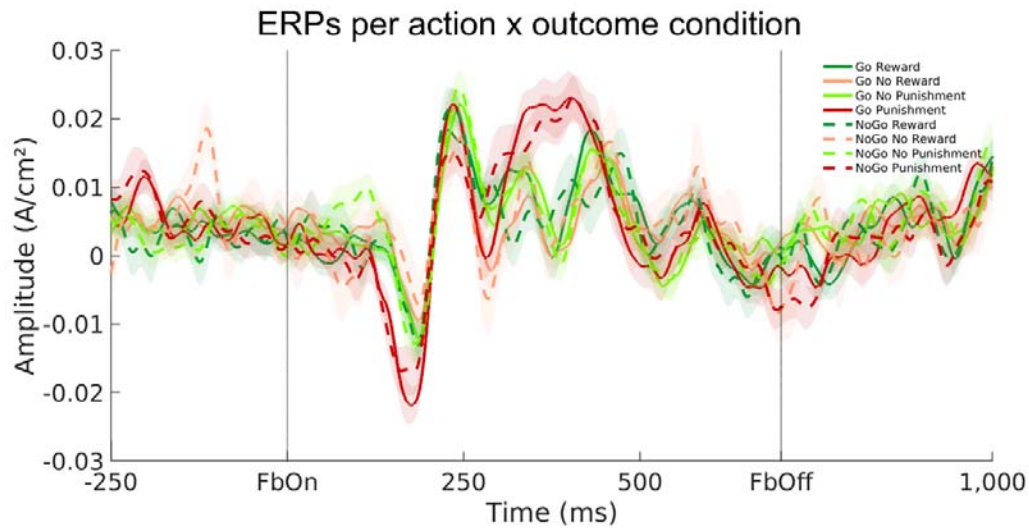
1688

**Figure S09B. ERPs per action x outcome condition**. Biggest differences occurred around the time of the N2 (FRN) and P3 (RewP). N2 and P3 exhibited larger amplitudes on trials with punishments. There was no apparent modulation by the previous action (Go/ NoGo).

1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713

## S10: Model-based EEG analyses in the time domain

In addition to testing whether midfrontal time-frequency power reflected signatures of biased learning (see main text), we also tested whether the midfrontal time domain signal reflected biased learning. Again, we used the standard PE term and the difference term to biased PEs as regressors in a multiple linear regression on each channel-time bin.

Focusing on midfrontal electrodes, and controlling for outcomes valence, first, the $PE_{STD}$ term was negatively correlated with midfrontal voltage around 529–575 ms (p = .039; Fig. S10B). Note that so late after outcome onset, signal was not part of any "classical" ERP component any more. Second, the $PE_{DIF}$ correlated negatively with midfrontal voltage around 123–166 ms (p = .029) in the time range of the N1 and later positively around 365–443 ms (p < .001; Fig S10C) in the time range of the P3/ RewP. Third, a similar pattern of correlations occurred for the $PE_{BIAS}$ term (Cluster 1: negative, 111–184 ms, p = .004; Cluster 2: positive, 346–449 ms, p < .001; Fig. S10A).  Fourth, around these same time windows, midfrontal voltage also encoded outcome valence itself, but with opposite sign (Cluster 1: positive, 99–184 ms, p < .001; Cluster 2: negative, 308–448 ms, p < .001; see S09).

In sum, similar to analyses of midfrontal power reported in the main text, PE sign and magnitude were encoded in midfrontal voltage around the same time, but with opposite polarity: Signal around the time of the N1 encoded PE sign positively, but PE magnitude negatively. Vice versa, signal around the time of the P3/ RewP encoded PE sign negatively, but PE magnitude positively. The same phenomenon of separate valence and magnitude encoding in midfrontal EEG signal has been reported before (28–30). Notably, magnitude encoding in midfrontal voltage emerged for the $PE_{BIAS}$ term, but not the $PE_{STD}$, indicating that this correlation was driven by the $PE_{DIF}$ term and that biased learning described midfrontal voltage better than standard learning. These results complement our findings of theta/delta power encoding outcome valence and magnitude with opposite polarities (see main text).
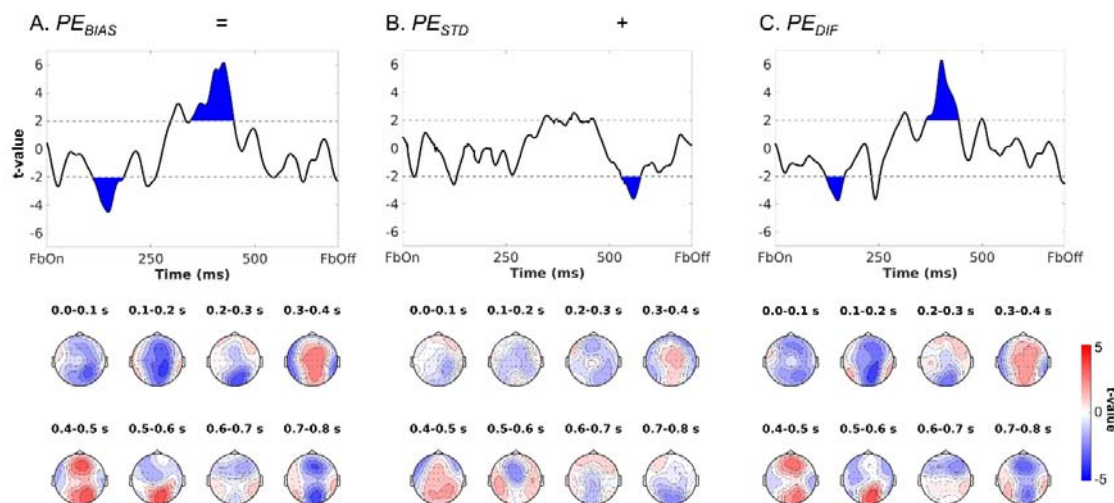


**Figure S10. Modulation of EEG voltage by biased PEs and decomposition into the standard PE term and the difference term to biased PEs.** (A) Mean EEG voltage over midfrontal electrodes (Fz, FCz, Cz) was significantly modulated by biased PEs around 111–184 (negatively) and 353–414 ms (positively) after outcome onset. (B) Correlations with the standard PE term only emerged around 529 – 575 ms (negatively). (C) Correlations with the difference term to biased PEs were similar to correlations for the biased PE term itself, i.e., around 123–166 (negatively) and 365–443 ms (positively).

Bottom row: Topoplots displaying *t*-values of beta-weights for the respective regressor over the entire scalp in steps of 100 ms from 0 to 800 ms.

## S11: Supplementary fMRI-inspired EEG results in time-frequency space

Besides the results for striatum, ACC, and PCC reported in the main text, there were also significant EEG correlates over midfrontal electrodes for trial-by-trial BOLD signal from left motor cortex ($p$ = .002, around 0–625 ms, 16–27 Hz; Fig. S11A). There were however no significant EEG correlates over midfrontal electrodes for BOLD signal from vmPFC/ subgenual ACC ($p$ = .174; Fig. S11B), left inferior temporal gyrus ($p$ = .097; Fig. S11C), and primary visual cortex ($p$ = .017; Fig. S11D).

As quality checks, we checked whether visual cortex BOLD correlated negatively with alpha over occipital electrodes (90, 91) and whether motor cortex BOLD correlated negatively with beta power over central electrodes (92, 93). Both was the case (see Fig. S11E and F), showing that our data was of sufficient quality to detect these well-established associations.
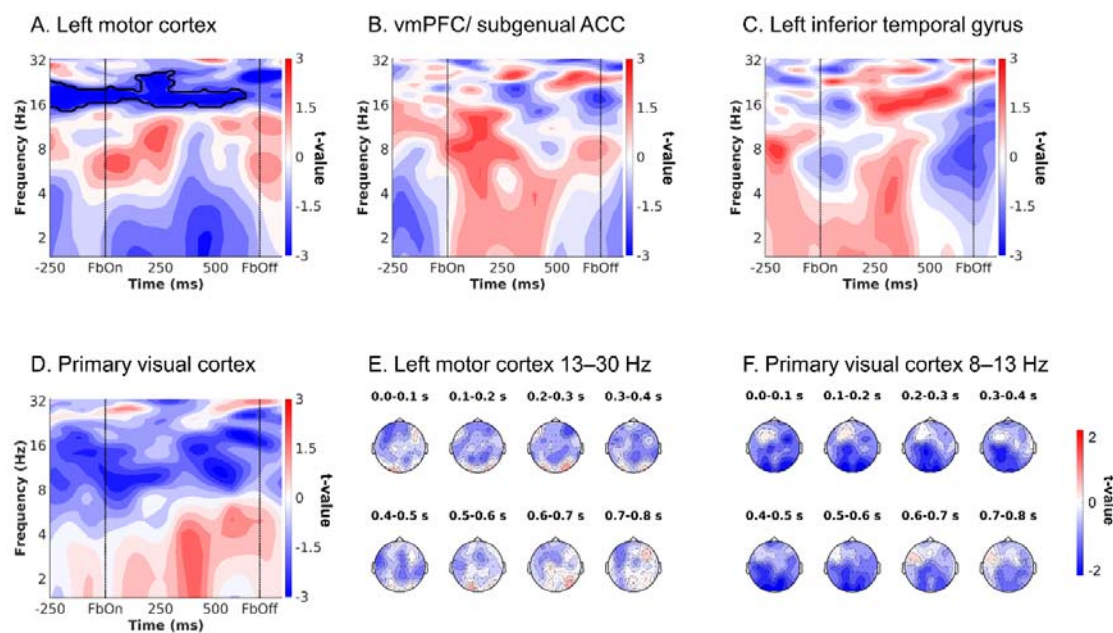


**Figure S11. Supplementary fMRI-informed EEG results in the time-frequency domain.** Unique temporal contributions of BOLD signal in (A) left motor cortex, (B) vmPFC, (C) left ITG and (D) primary visual cortex to midfrontal EEG power. Group-level $t$-maps display the modulation of the EEG power over midfrontal electrodes (Fz/ FCz/ Cz) by trial-by-trial BOLD signal in the selected ROIs. There significant correlations between midfrontal EEG TF power in the beta range and left motor cortex BOLD signal ($p$ = .002), but no significant midfrontal EEG correlates for BOLD signal from other ROIs. (E) Topoplot displaying $t$-values of left motor cortex BOLD over the entire scalp between 13 and 30 Hz (beta band) in steps of 100 ms from 0 to 800 ms. There are significant negatively correlates over central electrodes, especially round 300–500 ms. (F) Topoplot displaying $t$-values of primary visual cortex BOLD over the entire scalp between 8 and 13 Hz (alpha band) in steps of 100 ms from 0 to 800 ms. There are significant negatively correlates over occipital electrodes throughout outcome presentation.

## S12: Supplementary fMRI-inspired EEG results in the time domain

For fMRI-inspired analysis of the EEG signal in the time domain (voltage), we applied the same approach as reported in main text, but with voltage signal (time-domain) instead of time-frequency power as dependent variable. As independent variables, we entered the trial-by-trial BOLD signal from all seven regions encoding biased PEs plus the trial-by-trial standard PE and the different term towards the biased PE (exact same procedure as for EEG TF analyses), all in one single multiple linear regression. On a group-level, we again focused on the mean signal over midfrontal electrodes (Fz/ FCz/ Cz) in a time range of 0–700 ms, for which ERPs had been visible in the condition-averaged plots (see S09).

First, trial-by-trial striatal BOLD correlated significantly with midfrontal voltage at two time points, namely positively around 152–196 ms ($p$ = .017) in the time range of the N1 and again negatively around 316–383 ms ($p$ < .001, see Fig. S12A) in the time range of the N2/ FRN and P3/RewP. Second, trial-by-trial vmPFC BOLD correlated significantly positively with midfrontal voltage around 347–412 ms ($p$ = .006, see Fig. S12A) in the time range of the N2/ FRN and P3/RewP. Third, trial-by-trial BOLD from primary visual cortex correlated significantly positively with midfrontal voltage around 307–367 ms ($p$ = .011, see Fig. S12B), overlapping with (but slightly earlier than) correlations from vmPFC BOLD, i.e., in the time range of the N2/ FRN and P3/RewP. For midfrontal voltage split up per high vs. low BOLD signal (revealing which ERP components are respectively modulated), see Fig. S12C-E. There were no significantly correlations between midfrontal voltage and trial-by-trial BOLD from ACC ($p$ = .927, see Fig. S12A), left motor cortex ($p$ = .649, see Fig. S12B), PCC ($p$ = .796, see Fig. S12A), or left inferior temporal gyrus ($p$ = .649, see Fig. S12B). For further details on BOLD-EEG voltage correlations in the time domain, see Fig. S12F–L.

Taken together, trial-by-trial BOLD signal in striatum, vmPFC, and V1 all correlated with FRN/ RewP amplitude, which is the dominant phenomenon over midfrontal electrodes reflecting outcome valence (see S09 and S10). Notably, correlations with striatal and vmPFC BOLD were of opposite signs, which aligns with the finding that striatal and vmPFC BOLD predicted opposite behavioral tendencies on future trials (see main text; see S15). However, crucially, the time domain signal did not allow for a temporal dissociation of these different regions. Possibly, the midfrontal evoked signal (i.e., the part of the signal that is phase-locked to outcome onset) is so stereotyped that only the FRN/RewP complex shows enough variation across trials to allow for substantial correlations with trial-by-trial BOLD signal. This finding demonstrates that the time-frequency domain signal (i.e., the part of the signal that is not necessarily phase-locked to outcome onset) might be more suited for dissociating the activity of different regions in time.
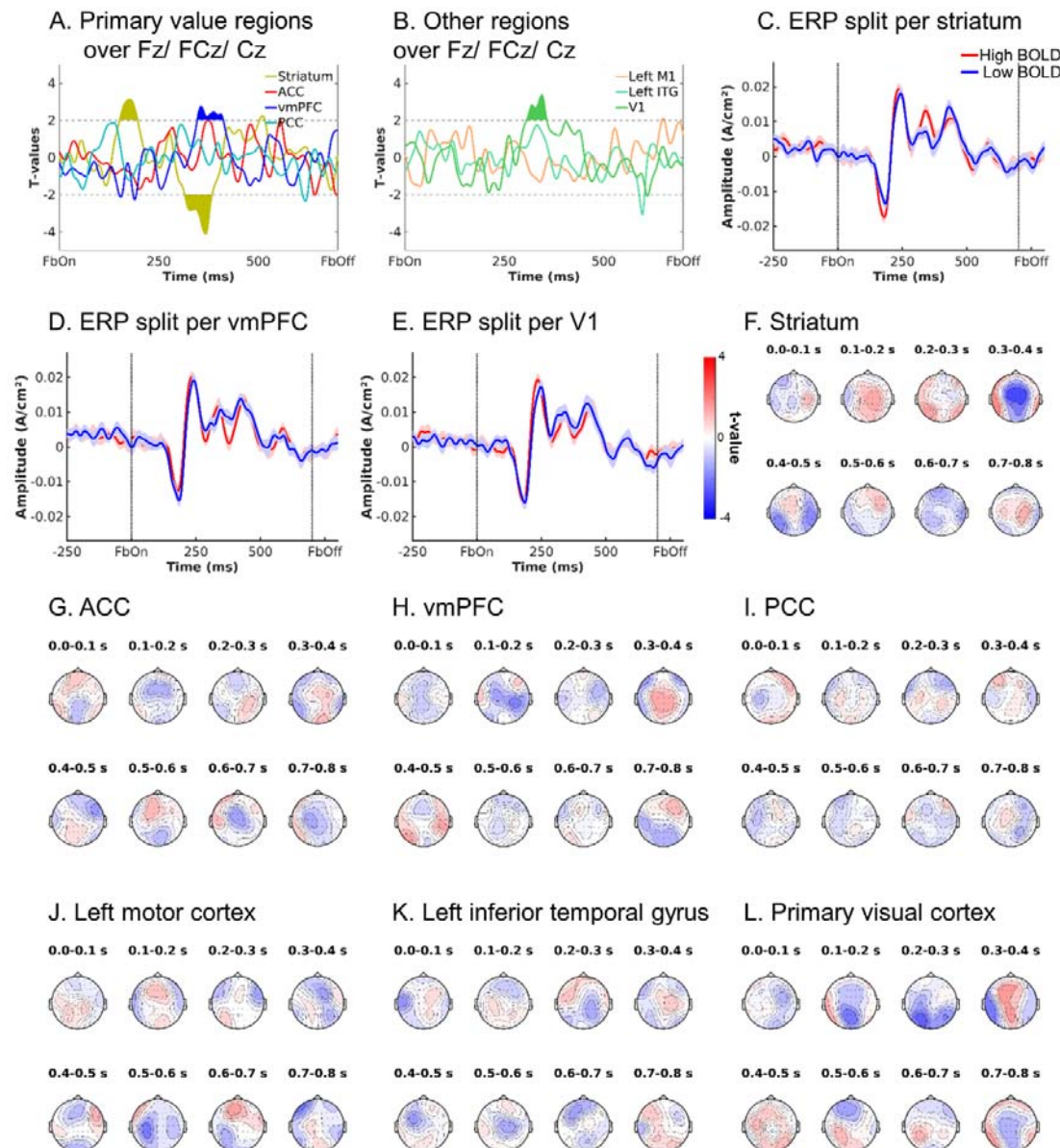
**Figure S12. fMRI-informed EEG analyses in the time-domain.** Group-level $t$-value time courses display the modulation of the EEG voltage over midfrontal electrodes (Fz/ FCz/ Cz) by trial-by-trial BOLD signal in the selected ROIs. (A) Correlations between midfrontal voltage and trial-by-trial BOLD signal from core value regions, i.e., striatum, ACC, vmPFC, and PCC. Striatal BOLD modulates the amplitude of the N1 and P3, while the P3 amplitude is also modulated by vmPFC BOLD. (B) Correlations between midfrontal voltage and trial-by-trial BOLD signal from other regions, i.e., left motor cortex, left inferior temporal gyrus, and primary visual cortex. Visual cortex BOLD modulates the amplitude of the P3, as well. (E-F) Midfrontal voltage split up for high vs. low BOLD signal (median split) from regions significantly modulating voltage. Striatal BOLD modulates N1 and P2 amplitude, while vmPFC BOLD and visual cortex BOLD modulate N2 (FRN) amplitude. (G-L) Topoplots displaying $t$-values of correlations between midfrontal voltage and trial-by-trial BOLD for all regions in steps of 100 ms from 0 to 800 ms.

1789

1790

1791

1792

1793 ## S13: Full list of significant clusters with EEG regressors in fMRI GLMs

| No | Contrast / Brain region | Maximal Z-value | Cluster size (voxels) | Corrected p | Peak coordinates x | y | z |
|----|------|------|------|------|------|------|------|
| | **Central Lower Alpha Band Positive** | | | | | | |
| | *No significant clusters* | | | | | | |
| | **Central Lower Alpha Band Negative** | | | | | | |
| 1 | Precuneous, cuneal cortex, right superior lateral occipital cortex | 5.78 | 8346 | 2.50e-33 | 6 | -60 | 66 |
| 2 | Anterior cingulate gyrus, right superior frontal gyrus | 4.77 | 2449 | 1.75e-14 | 24 | 12 | 66 |
| 3 | Left middle frontal gyrus, | 5.59 | 1828 | 7.63e-12 | -38 | 8 | 34 |
| 4 | Right insula, right central opercular cortex | 4.71 | 1794 | 1.08e-11 | 42 | 2 | 28 |
| 5 | Right frontal pole, right middle frontal gyrus, right inferior frontal gyrus, pars triangularis | 5.43 | 1300 | 2.37e-09 | 30 | 40 | 20 |
| 6 | Left supramarginal gyrus, anterior division | 4.61 | 959 | 1.19e-07 | -64 | -36 | 42 |
| 7 | Left angular gyrus | 5.83 | 916 | 2.38e-07 | -48 | -52 | 18 |
| 8 | Right cerebellum, anterior | 4.79 | 480 | .000131 | 42 | -38 | -38 |
| 9 | Posterior cingulate cortex, parahippocampal gyrus, right thalamus | 4.41 | 424 | .000328 | 14 | -38 | -2 |
| 10 | Left temporal pole, left inferior frontal gyrus, pars opercularis left insula | 4.08 | 413 | .000394 | -56 | 16 | -6 |
| 11 | Left cerebellum, anterior | 5.44 | 263 | .00598 | -30 | -40 | -42 |
| 12 | Right lingual gyrus | 3.43 | 235 | .0104 | 10 | -74 | -10 |
| 13 | Left cerebellum, posterior | 5.74 | 215 | .0158 | -14 | -76 | -42 |
| 14 | Brainstem | 4.35 | 207 | .0186 | 8 | -34 | -20 |
| | **Frontal Theta Band Positive** | | | | | | |
| 1 | Right bilateral precentral gyrus | 4.82 | 394 | .000577 | 12 | -16 | 80 |
| 2 | Left bilateral precentral gyrus | 5.25 | 357 | .0011 | -20 | -28 | 78 |
| | **Frontal Theta Band Negative** | | | | | | |
| 1 | Right supramarginal gyrus, posterior division, right superior lateral occipital cortex | 3.94 | 1002 | 1.10e-07 | -54 | -50 | 44 |
| 2 | Left supramarginal gyrus, posterior division, Left superior lateral occipital cortex | 4.39 | 508 | 8.96e-05 | 56 | -50 | 20 |
| 3 | Posterior cingulate cortex | 4.58 | 419 | .000378 | -6 | -30 | 38 |
| 4 | Ventromedial prefrontal cortex | 4.03 | 342 | .00143 | 0 | 42 | 4 |
| | **Central Beta Band Positive** | | | | | | |
| 1 | Right caudate | 4.19 | 258 | .00481 | 16 | 30 | 6 |
| 2 | Left parahippocampal gyrus, posterior divison | 4.86 | 221 | .0106 | -38 | -36 | -8 |
| | **Central Beta Band Negative** | | | | | | |
| 1 | Right frontal pole, right middle frontal gyrus, right superior frontal gyrus | 5.49 | 6599 | 7.06e-30 | -32 | 8 | 28 |
| 2 | Left frontal pole, left middle frontal gyrus, Left superior frontal gyrus | 5.51 | 6144 | 1.82e-28 | 40 | 38 | 36 |
| 3 | Left supramarginal gyrus, posterior | 5.51 | 5175 | 2.43e-25 | -66 | -44 | 28 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | division, left superior parietal lobule, left superior lateral occipital cortex, Left middle temporal gyrus, temporooccipital part | | | | | | |
| 4 | Right supramarginal gyrus, posterior division, Right superior parietal lobule, right superior lateral occipital cortex | 5.13 | 3264 | 1.62e-18 | 30 | -74 | 54 |
| 5 | Left superior frontal gyrus, paracingulate gyrus, precuneous | 4.54 | 1235 | 1.80e-09 | -4 | 12 | 52 |
| 6 | Right superior temporal gyrus, posterior division | 4.59 | 1076 | 1.33e-08 | 48 | -14 | -10 |
| 7 | Left temporal pole, left planum temporale | 4.96 | 320 | .00139 | -46 | 4 | -18 |

1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823

## 1824    S14: Go/NoGo difference in alpha (and beta) over time

1825    We observed differences between trials with Go responses and trials with NoGo responses in
1826    the low alpha power before and shortly after outcome onset (Fig. 6A, B main text). Alpha typically

1827   increases over the time course of an experiment, potentially related to fatigue and decreasing
1828   arousal (*94*). If the ratio of Go and NoGo responses changed over time, as well, such an increase over
1829   time could spuriously lead to a difference between Go and NoGo responses (though note that this
1830   ratio did not noticeably change over time; Fig. S14D). To exclude this possibility, we extracted trial-
1831   by-trial time-frequency power from the three significant clusters report in the main text in which
1832   power differed between Go and NoGo responses: a) lower alpha band power after outcome onset, b)
1833   lower alpha band power before and after outcome onset, c) beta band power before outcome onset.
1834   We transformed this data to decibel and analyzed it as a function of the performed response (factor),
1835   block number (1–6; z-standardized), and the interaction between both. We reasoned that if power
1836   differences occurred merely due to fatigue effects, the main effect of performed response should not
1837   be significant when accounting for time on task (i.e., block number).
1838         For lower alpha band power after outcome onset, there was a significant main effect of
1839   performed response, $b = 0.035$, $SE = 0.015$, $\chi^2(1) = 5.350$, $p = .021$, with higher power for Go than
1840   NoGo responses, a significant main effect of block number with lower alpha band power increasing
1841   over time, $b = 0.052$, $SE = 0.019$, $\chi^2(1) = 6.645$, $p = .010$, but no significant interaction, $b = 0.003$, $SE =$
1842   $0.008$, $\chi^2(1) = 0.156$, $p = .693$. As Fig. S14A reveals, lower alpha band power was consistently higher
1843   after Go than after NoGo responses for every block of the task, suggesting that differences in lower
1844   alpha band power were not merely due to time on task.
1845         For lower alpha band power before and after outcome onset, as well, there was a significant
1846   main effect of performed response, $b = 0.068$, $SE = 0.030$, $\chi^2(1) = 5.010$, $p = .025$, with higher power
1847   after Go than NoGo responses, a significant main effect of block number with lower alpha band
1848   power increasing over time, $b = 0.072$, $SE = 0.029$, $\chi^2(1) = 6.757$, $p = .016$, but no significant
1849   interaction, $b = 0.010$, $SE = 0.009$, $\chi^2(1) = 1.184$, $p = .277$ (Fig. S14B), leading to identical conclusions.
1850         For beta band power before and after outcome onset, there was a significant main effect of
1851   performed response, $b = 0.083$, $SE = 0.032$, $\chi^2(1) = 6.301$, $p = .012$, with higher power after Go than
1852   NoGo responses, a significant main effect of block number with beta power decreasing over time, $b =$
1853   $-0.042$, $SE = 0.021$, $\chi^2(1) = 4.007$, $p = .045$, but no significant interaction, $b = 0.001$, $SE = 0.007$, $\chi^2(1) =$
1854   $0.030$, $p = .864$ (Fig. S14C). In sum, even in presence of changes in power over the time course of the
1855   task, lower alpha band and beta band power were consistently higher after Go responses than after
1856   NoGo responses, suggesting that these effects were not due to time on task.
1857         Furthermore, we asked whether differences in ACC BOLD between trials with Go and trials
1858   with NoGo response at the time of the outcome were due to outcome-related activity or might
1859   rather the reflect action on the next trial. We thus plotted the "raw" BOLD signal per action x
1860   outcome condition. We used the first eigenvariate of the BOLD in signal in the ACC cluster that
1861   reflected biased learning, upsampled the BOLD signal, epoched it into trials relative to outcome
1862   onset (same procedure as for fMRI-informed EEG analyses), and averaged the signal across trials and
1863   participants separately per performed action (Go/NoGo) and outcome valence (positive/ negative).
1864   This plot yielded higher ACC BOLD signal on trials with NoGo responses than on trials with Go
1865   responses at the time of outcomes (Fig. S14E). However, this difference could potentially be driven
1866   by the response on the following task, so we further split the data according to whether the action
1867   on the following trial was a Go or a NoGo response. Irrespective of the action on the following trial,
1868   ACC BOLD signal was higher when the action on the current trial was a NoGo response compared to a
1869   Go response (Fig. S15F). In sum, these analyses corroborate that ACC BOLD signal was indeed higher
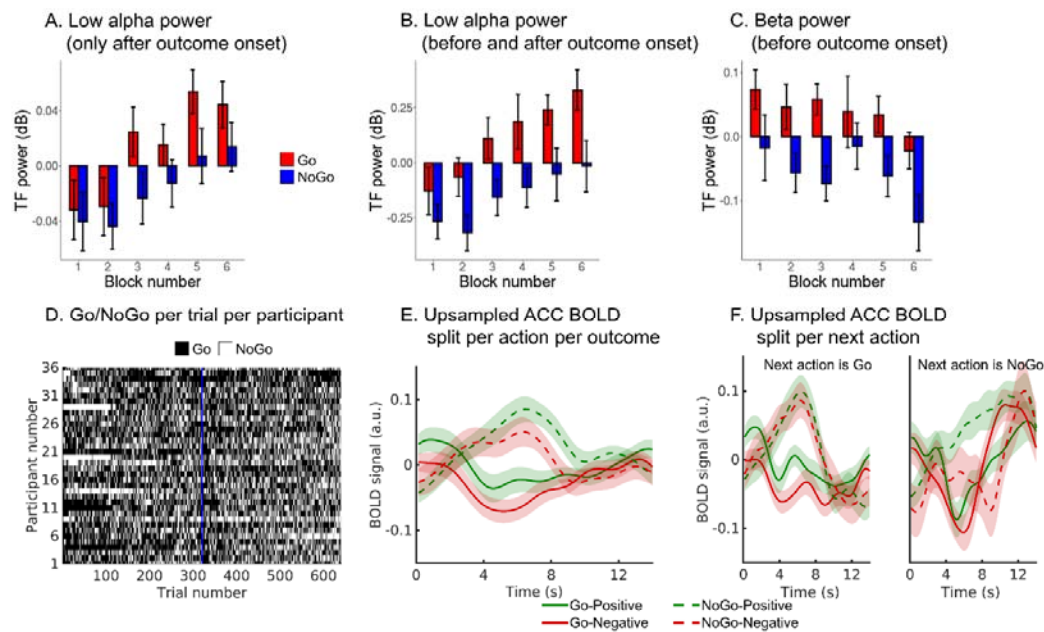1870   after NoGo than Go responses at the time of outcomes.
1871

**Figure S14. Control analyses excluding temporal confounds in midfrontal lower alpha band power and ACC BOLD**. (A) Mean midfrontal low alpha power (±SEM across participants) after outcome onset, (B) before and after outcome onset, and (C) beta power before outcome onset as a function of the performed action and block number (i.e., time on task). While low alpha power increases and beta power decreases over the time course of the task, power is always consistently higher for trials with Go than trials with NoGo responses, suggesting that action effects are not reducible to time on task. (D) Response for each participant (rows) on each trial (columns). There is no noticeable change in the overall ratio of Go to NoGo responses over time. The vertical blue line indicates the start of the second session featuring new stimuli. (E) Mean upsampled ACC BOLD signal (±SEM across participants) at the time of the outcome, split per performed action (Go/NoGo) and outcome valence (positive/negative). BOLD signal is higher after NoGo than Go responses. (F) Same plot as (E), but split based on whether the next action is a Go (left panel) or an NoGo (right panel) response. Even if the next response is NoGo, BOLD signal is higher for trials with NoGo responses (on the current trial) than trials Go responses.

1872

1873

1874

1875    ## S15: Stay behavior as a function of BOLD and EEG TF power
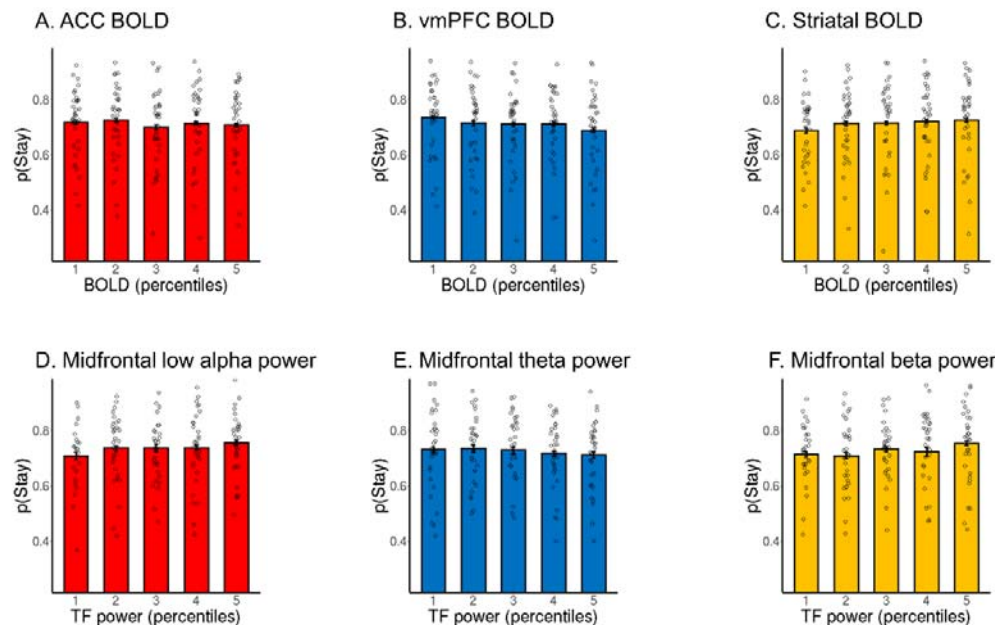


**Figure S15. Probability of repeating the same response ("stay") on the next cue encounter as a function of outcome-related BOLD and EEG signal.** (A-C) Probability of repeating the same action ("staying") as a function of BOLD signal from (A) ACC, (B) vmPFC, and (C) striatum (split into 5 bins). While ACC BOLD was not significantly linked to the probability to stay, high BOLD signal in vmPFC predicted a higher chance to switch to another action, while high BOLD signal in striatum predicted a higher probability of staying with the same action. (D-E) Probability of staying as a function of midfrontal time-frequency power in the (A) low alpha, (B) theta/delta, and (C) beta range. Higher low alpha power and higher beta power predict a higher probability of staying with the same action, while higher theta power predicts a higher chance to switch to another action. Grey circles represent individual per condition-per-participant means. Error bars are very narrow (and thus hardly visible) and computed based on the Cousineau-Morey methods based on per-condition-per-participant means.

1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889