

1 **TITLE**

2 Valence-partitioned learning signals drive choice behavior and phenomenal subjective experience in  
3 humans

4  
5  
6 **AUTHORS**

7 L. Paul Sands<sup>1,2,\$</sup>, Angela Jiang<sup>1</sup>, Rachel E. Jones<sup>1,2</sup>, Jonathan D. Tractner<sup>1,2</sup>, Kenneth T. Kishida<sup>1,2,3\*</sup>

8  
9 **AFFILIATIONS**

10 <sup>1</sup>Dept. of Physiology and Pharmacology, Wake Forest School of Medicine, Winston-Salem NC, 27101,  
11 US

12 <sup>2</sup>Neuroscience Graduate Program, Wake Forest School of Medicine, Winston-Salem NC, 27101, US

13 <sup>3</sup>Dept. of Neurosurgery, Wake Forest School of Medicine, Winston-Salem NC, 27101, US

14  
15 **\*CORRESPONDING AUTHOR**

16 Kenneth T. Kishida, PhD; Wake Forest School of Medicine, Medical Center Blvd, Winston-Salem NC,  
17 27101, US. phone: 336-716-0419, email: [kkishida@wakehealth.edu](mailto:kkishida@wakehealth.edu)

18  
19 <sup>\$</sup>**Current address:** Fralin Biomedical Research Institute at Virginia Tech

20  
21 **Highlights:**

- 22
- 23 • TD-Reinforcement Learning (RL) theory interprets punishments relative to rewards.
  - 24 • Environmentally, appetitive and aversive events are statistically independent.
  - 25 • Valence-partitioned RL (VPRL) processes reward and punishment independently.
  - 26 • We show VPRL better accounts for human choice behavior and associated BOLD activity.
  - 27 • VPRL signals predict dynamic changes in human subjective experience.

28 **SUMMARY**

29 How the human brain generates conscious phenomenal experience is a fundamental problem. In  
30 particular, it is unknown how variable and dynamic changes in subjective affect are driven by interactions  
31 with objective phenomena. We hypothesize a neurocomputational mechanism that generates valence-  
32 specific learning signals associated with ‘what it is like’ to be rewarded or punished. Our hypothesized  
33 model maintains a partition between appetitive and aversive information while generating independent  
34 and parallel reward and punishment learning signals. This valence-partitioned reinforcement learning  
35 (VPRL) model and its associated learning signals are shown to predict dynamic changes in 1) human  
36 choice behavior, 2) phenomenal subjective experience, and 3) BOLD-imaging responses that implicate  
37 a network of regions that process appetitive and aversive information that converge on the ventral  
38 striatum and ventromedial prefrontal cortex during moments of introspection. Our results demonstrate  
39 the utility of valence-partitioned reinforcement learning as a neurocomputational basis for investigating  
40 mechanisms that may drive conscious experience.

41  
42 **KEYWORDS:** consciousness, subjective experience, decision-making, reinforcement learning, reward  
43 prediction errors, punishment prediction errors, valence, affect

## 44 INTRODUCTION

45  
46 The mechanisms by which the human brain generates the subjective phenomenal experiences  
47 that allow us to answer introspective questions like, “What is it like to be [me]?” (or “a bat”; Nagel, 1974)  
48 remain a fundamental mystery that has occupied artists, philosophers, and neuroscientists for centuries  
49 (Faherty, 2016). However, this problem represents more than just an old philosophical quandary: brain  
50 states underlying subjective suffering and challenges to the ability to control one’s behavior are at the  
51 core of nearly all psychiatric and neurological conditions (Kishida et al., 2010; Montague et al., 2012;  
52 Kishida, 2012; Redish and Gordon, 2016; Huys et al., 2016; Taschereau-Dumouchel et al., 2022). The  
53 inherently subjective nature of conscious experience has led philosophers to deem an understanding of  
54 the mechanisms supporting it fundamentally ‘hard’ or even impossible (Nagel, 1974; Chalmers, 1995).  
55 On the other hand, empirical investigation has turned previously seemingly impossible problems (e.g.,  
56 an understanding of electromagnetic phenomena; Forbes and Mahon, 2014) into well-defined scientific  
57 fields of inquiry. Here, “we get on with the task” of empirically investigating simple conscious experiences  
58 through the lens of behavioral and neurobiological measurements (Churchland PM, 1984, 2014;  
59 Churchland PS, 1996) that may be better understood within a neurocomputational framework  
60 (Churchland and Sejnowski, 1994; Kishida, 2012). One of the major challenges facing a science of  
61 consciousness lies in the fact that the phenomena to be investigated – e.g., variations in how one feels  
62 – are subjective and only *indirectly* accessible through self-report behavior. Nonetheless, subjective  
63 experiences are associated with reproducible behaviors and changes in neurophysiology that can be  
64 studied within behavioral, cognitive, and computational neuroscience methods.  
65

66 A leading neurocomputational approach to investigating adaptive human choice behaviors and  
67 subjective experiences has been the use of temporal difference (TD) reinforcement learning (RL) theory  
68 (Sutton, 1988; Sutton and Barto, 1998) to provide a framework for probing how dopamine neurons  
69 encode ‘teaching signals’ in the form of TD reward prediction errors (RPEs; Montague et al., 1996; Schultz  
70 et al., 1997; Bayer and Glimcher, 2005; Bayer et al., 2007; Zaghoul et al., 2009; Glimcher, 2011; Hart  
71 et al., 2014; Eshel et al., 2015; Kishida et al., 2016; Watanabe-Uchida et al., 2017; Moran et al., 2018).  
72 In the dopamine TD-RPE hypothesis (Montague et al., 1996; Schultz et al., 1997), phasic bursts and  
73 pauses in dopamine neuron firing activity signal ‘better-than-expected’ or ‘worse-than-expected’  
74 prediction errors, respectively, which provides a computationally-optimal method for learning – directly  
75 from experience – value associations between rewards and the stimuli and actions that predict them  
76 (Sutton and Barto, 1998). This mechanistic insight has since led to specific hypotheses about the  
77 neurochemical basis of computations underlying human choice behaviors and a variety of mental health  
78 disorders (Redish and Gordon, 2016), in part due to support from human fMRI studies demonstrating  
79 that BOLD activity in brain regions rich in dopaminergic terminals parametrically tracks reward prediction  
80 errors during classical and instrumental conditioning (O’Doherty et al., 2003; McClure et al., 2003;  
81 Pessiglione et al., 2006; Garrison et al., 2013). Furthermore, empirical studies have begun to associate  
82 neurocomputational processes underlying RPE encoding with the immediate subjective experience of  
83 pleasure as well as associated dynamic changes in mood that occur over longer timescales (Delgado et  
84 al., 2006; Xiang et al., 2013; Rutledge et al., 2014; Eldar et al., 2016). This work has also provided a  
85 basis for investigating the neural and behavioral correlates of changes associated with various psychiatric  
86 conditions and mood disorders (Redish and Gordon, 2016; Redish, 2004; Montague et al., 2012; Kishida  
87 et al., 2010; Huys et al., 2016; Rutledge et al., 2017; Brown et al., 2021).  
88

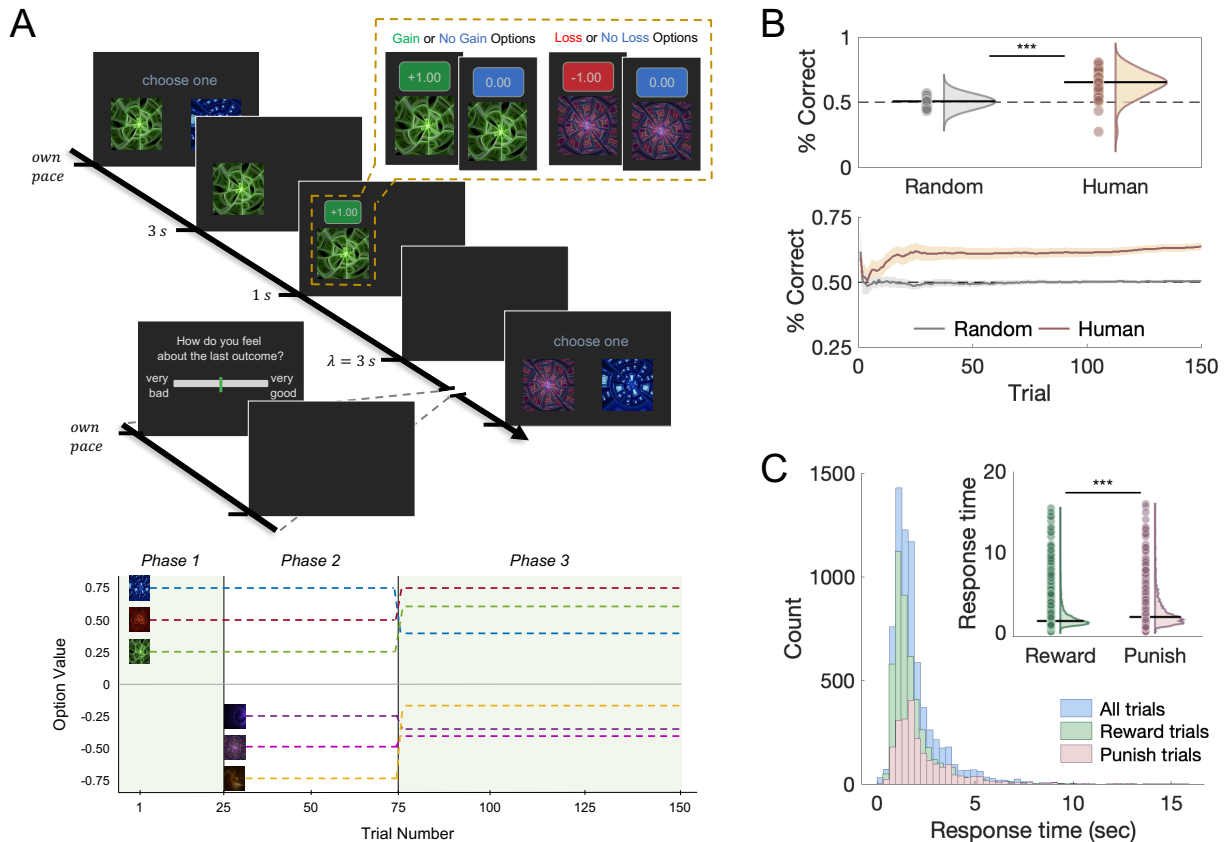
89 Despite its overwhelming utility, RL theory does not explicitly describe how biological organisms  
90 learn optimally from aversive experiences (i.e., punishments) concurrently or independently from  
91 experienced rewards (Sutton, 1988; Sutton and Barto, 1998; Dayan and Niv, 2008; Pessiglione and  
92 Delgado, 2015; Kishida and Sands, 2021). Aversive stimuli (e.g. those that cause tissue damage or  
93 threaten to do so) are evolutionarily conserved drivers of defensive and other adaptive behaviors (Cisek,  
94 2021) and negative aspects of human subjective experience (Seymour et al., 2007b; Kishida and Sands,  
95 2021). Typically, TDRL theory treats aversive experiences (e.g., punishments or costs) as ‘negative  
96 rewards’ and thus colinear with rewards along a single valence dimension: TDRL treats punishing (i.e.,  
97 aversive) outcomes only *in relation to* appetitive or rewarding experiences. This is counter to biological

98 experience where appetitive and aversive experiences may largely be derived from statistically  
99 independent sources or derived from a similar source with varying degrees of statistical dependence. For  
100 example, individual people may be collaborative or fiercely competitive; potential sources of food (e.g.,  
101 plants) may be nutritious or toxic. Further, computationally, using a single valence dimension ordains a  
102 'zero-sum rule' for how to represent co-occurring positive and negative outcomes (or mixed-valence  
103 experiences) – only a resultant single scalar 'reward' term is used in TDRL (Sutton and Barto, 1998) –  
104 and thereby also does not allow for dissociating the individual effects of co-occurring rewards and  
105 punishments on agent behavior and subjective experiences. Indeed, this traditional reliance on a  
106 unidimensional, colinear valence representation belies the independent influences of rewards and  
107 punishments (or otherwise appetitive and aversive events) on human choices and emotional responses  
108 (Konorski, 1967; Dickinson and Dearing, 1979; Cacioppo et al., 1999; Folkman and Moskowitz, 2000;  
109 Larsen et al., 2009; Larsen and McGraw, 2011; Kishida and Sands, 2021).

110  
111 Fundamentally, there remains a gap in the literature regarding traditional TDRL accounts of  
112 reward and punishment learning and comparisons to alternative models that directly address how  
113 punishing stimuli may be processed in a comparably optimal manner. Human fMRI investigations of  
114 aversive valence-processing suggest that adaptive learning from punishments (e.g., pain) is consistent  
115 with a hypothetical punishment-based (i.e., reward-opponent) RL system (Palminteri and Pessiglione,  
116 2017; Seymour et al., 2004, 2005, 2007a, 2012; Delgado, et al. 2008, 2009). Theoretical descriptions of  
117 a system 'opponent' to dopaminergic reward processing have been hypothesized (Daw et al., 2002) and  
118 are supported by indirect evidence (Palminteri and Pessiglione, 2017; Seymour et al., 2005, 2007a,  
119 2007b; Delgado, et al. 2008, 2009) and recent direct simultaneous measurements of serotonin and  
120 dopamine in human striatum (Kishida et al., 2016; Moran et al., 2018). However, these prior investigations  
121 generally used a unidimensional representation of valence. To begin to explicitly compare alternatives to  
122 unidimensional TDRL-based depictions of reward and punishment learning, we hypothesized *valence-*  
123 *partitioned reinforcement learning* (VPRL; Kishida and Sands, 2021), which proposes that separate  
124 neural systems implement TD learning over appetitive and aversive experiences in parallel and thereby  
125 independently update representations of positive and negative expected state-action values,  
126 respectively. VPRL-encoded signals can then be operated on (e.g., integrated) or processed  
127 independently as necessary for guiding behavior, including when introspecting or reporting about one's  
128 subjective feelings (Kishida and Sands, 2021).

129  
130 Here, we test the hypothesis that VPRL is a better model than traditional TDRL for investigating  
131 (1) human learning and decision-making behavior, (2) associated neural activity, and (3) dynamic  
132 moment-to-moment changes in subjective experience in humans. We combine data from two  
133 experiments involving human participants (N=47 total) scanned with fMRI while performing a probabilistic  
134 reward and punishment (PRP) task that uses monetary gains and losses as reinforcement (**Figure 1A**;  
135 Methods). We show that VPRL better explains participant choice behavior compared to traditional TDRL  
136 and that VPRL model parameters fit to participant choices are consistent with humans learning from  
137 rewards and punishments independently and asymmetrically (**Figure 2**). Further, we investigate the  
138 connection between VPRL learning signals and participants' self-reported subjective feelings about  
139 received outcomes throughout the PRP task, demonstrating that the expected value of a chosen action  
140 and prediction errors over action-contingent rewards and punishments all uniquely influence – and  
141 together predict – subjective feelings about experienced outcomes (**Figure 3**). Model-based fMRI  
142 analyses reveal blood-oxygen-level-dependent (BOLD) signals that parametrically track VPRL learning  
143 signals and associated subjective feelings within a distributed network of striatal, cingulate, insular, and  
144 prefrontal regions (**Figure 4,5**). Our results support the notion of valence partitioning as a mechanism in  
145 the human brain for processing appetitive and aversive stimuli via independent and parallel TDRL  
146 mechanisms, which together provide more robust representations of independent appetitive and aversive  
147 value estimates in uncertain contexts. Further, our results demonstrate and we discuss the implications  
148 of VPRL as a valid neurocomputational framework for investigating the neural mechanisms underlying  
149 the dynamics of subjective phenomenal experience and associated behaviors in humans.

150  
151



**Figure 1 – Human performance on a probabilistic reward and punishment (PRP) task.** (A, top) Schematic of a trial from the PRP task and subjective rating prompt. On each trial, a participant chooses one of two options and is reinforced probabilistically with either a monetary gain, nothing, or a monetary loss. Randomly after a third of trials, participants submit ratings of their subjective feelings about experienced outcomes. Offset: reward-associated options result in either monetary gains or nothing, and punishment-associated options result in monetary losses or nothing. (A, bottom) Depiction of the ‘ground-truth’ expected value for each option (expected value = probability(outcome)\*outcome) and how the options’ expected values change throughout the three phases of the PRP task (demarcated by vertical black lines). In phase 1 you choose between two of 3 possible gain/no-gain options. For phase 2, there’s an equal number of trials with two gain/no-gain options and two loss/no-loss. In phase 3, participants choose between any two of the six options at random, and the expected value for each option changes. Icon to outcome mappings are randomized for each participant. (B, top) The overall percent of trials where participants correctly chose the option with the highest (most positive) expected value, and (B, bottom) the evolution of the percent of correct choice trials throughout the PRP task. (C) The distributions of response times for all trials, trials on which a reward-associated option was chosen (reward trials), and trials on which a punishment-associated option was chosen (punish trials). For (B) and (C), \*\*\* =  $p < 0.001$ .

152

153

154

## RESULTS

155

156

### ***VPRL best explains choice behaviors and reveals asymmetrical processing of rewards and punishments***

157

158

159

160

161

162

Forty-seven participants completed the PRP task which required them to intermittently (randomly on one-third of all trials) rate how they felt about recent outcomes (Figure 1A). Participants learned to choose the option with the highest expected value on each trial more often than expected by chance (Figure 1B; two-sample t-test,  $t(92)=8.8$ ,  $p < 0.001$ ), chose the option with the highest expected value

163 increasingly over time (two-way ANOVA (group, time),  $F(1,149)=2416.8$ ), and were quicker to select  
 164 rewarding options than punishing options (**Figure 1C**; two-sample t-test,  $t(7023)=-14.5$ ,  $p<0.001$ ).

165  
 166 To test whether participants might learn differently from rewards and punishments, we fit a  
 167 standard TDRL model and a VPRL model (Kishida and Sands, 2021) to participant choice behavior using  
 168 hierarchical Bayesian inference and compared estimates of the model evidence (i.e., marginal likelihood)  
 169 and the posterior predictive accuracy (density) for both models. For both cohorts individually, and for an  
 170 'internal meta-analysis' combined cohort, VPRL demonstrated both greater model evidence and greater  
 171 posterior predictive accuracy relative to TDRL (**Table 1**), indicating that VPRL is a better explanation of  
 172 behavior on the PRP task and better predicts unobserved PRP task choice behavior data.

173  
 174

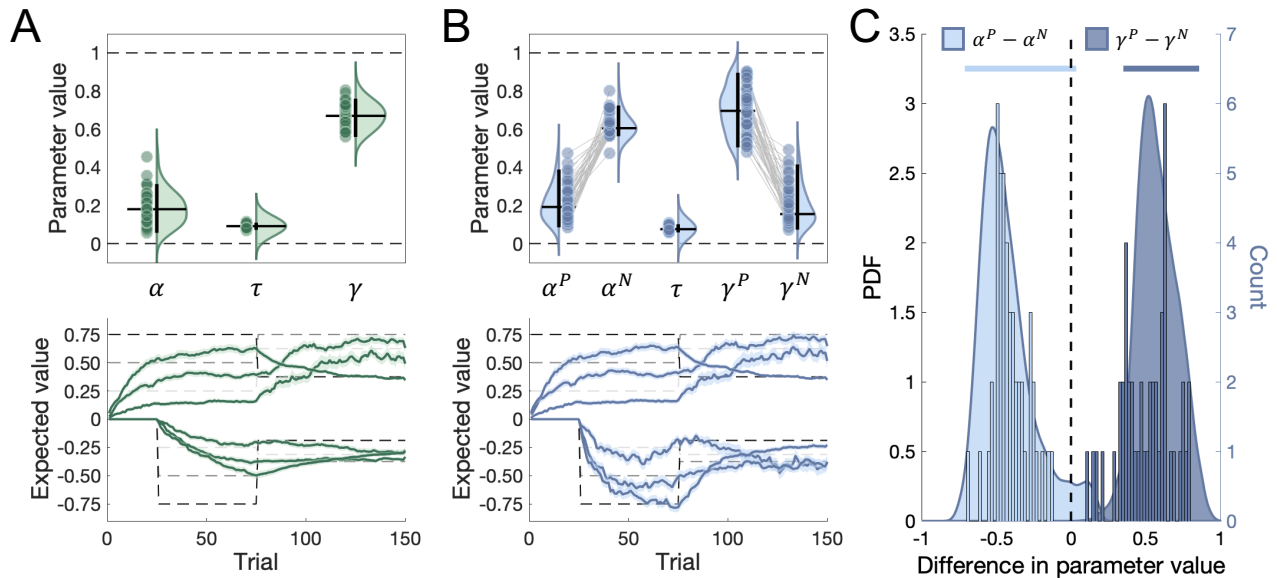
|            | fMRI Cohort 1 (n=20)            |                                 | fMRI Cohort 2 (n=27)            |                                 | Combined (N=47)                 |                                  |
|------------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|----------------------------------|
|            | <i>Model evidence</i>           | <i>Predictive density</i>       | <i>Model evidence</i>           | <i>Predictive density</i>       | <i>Model evidence</i>           | <i>Predictive density</i>        |
| TDRL       | -1577.3<br>(0.05)               | -1564.9<br>(56.8)               | -2400.2<br>(0.20)               | -2378.5<br>(61.7)               | -3979.3<br>(0.33)               | -3945.8<br>(89.9)                |
| VPRL       | <b>-1499.2</b><br><b>(0.24)</b> | <b>-1476.8</b><br><b>(63.3)</b> | <b>-2321.8</b><br><b>(0.48)</b> | <b>-2295.4</b><br><b>(71.0)</b> | <b>-3815.2</b><br><b>(2.45)</b> | <b>-3769.0</b><br><b>(101.9)</b> |
| Difference | -78.1<br>(0.25)                 | -88.2<br>(28.9)                 | -78.4<br>(0.35)                 | -83.0<br>(22.3)                 | -164.1<br>(2.25)                | -1767.7<br>(37.7)                |

**Table 1 – TDRL and VPRL model comparison results for two neuroimaging cohorts.** Computed estimates of the Bayesian model evidence (i.e., marginal likelihood) and model predictive density (i.e., cross-validated error) for TDRL and VPRL models. VPRL demonstrated the maximum (least negative) model evidence and expected predictive density (bold values) compared to TDRL. Reported values are the median estimate value (log scale), with values in parentheses indicating either the interquartile range (of model evidence estimations) or the Monte Carlo sampling error (for the predictive density). Given the similarity of TDRL and VPRL model comparison results for both fMRI cohorts separately and the improved model evidence and predictive density for VPRL when combining both fMRI cohorts, we elected to combine the data from both fMRI cohorts to improve the power of both the main behavioral and model-based fMRI analyses.

175

176 Given that participant choice behavior on the PRP task is better explained by VPRL compared to  
 177 TDRL, we next investigated differences in the (posterior) parameter distributions and the time series of  
 178 learned state-action values (Q-values) derived from each model (**Figure 2**). The group-level TDRL model  
 179 parameters are (**Figure 2A**): learning rate ( $\alpha$ ): median = 0.16 (95% credible (highest density) interval (CI)  
 180 = [0.13 0.21]); temporal discount factor ( $\gamma$ ): median = 0.65 [0.46 0.98]; and choice temperature ( $\tau$ ): median  
 181 = 0.09 [0.05 0.20]. The group-level VPRL model parameters are (**Figure 2B**): Positive valence (i.e.,  
 182 reward) system learning rate ( $\alpha^P$ ): median = 0.20 [0.16 0.24]; Negative valence (i.e., punishment) system  
 183 learning rate ( $\alpha^N$ ): median = 0.66 [0.17 0.91]; Positive system temporal discount factor ( $\gamma^P$ ): median =  
 184 0.71 [0.54 0.99]; Negative system temporal discount factor ( $\gamma^N$ ): median = 0.15 [0.03 0.29]; and choice  
 185 temperature ( $\tau$ ): median = 0.07 [0.05 0.14]. To investigate the nature of the differential learning of rewards  
 186 relative to punishments in the VPRL framework, we assessed the difference between the Positive and  
 187 Negative systems' learning rates and temporal discount parameters (**Figure 2C**). We found that the  
 188 learning rate for punishments is generally greater than the learning rate for rewards ( $\alpha^P - \alpha^N$  median  
 189 difference = -0.47 [-0.71 0.04]), and that temporal discounting for punishments was greater than temporal  
 190 discounting for rewards ( $\gamma^P - \gamma^N$  median difference = 0.56 [0.35 0.86]). Lastly, the time series of learned  
 191 expected values for TDRL (**Figure 2A**, bottom) and VPRL (**Figure 2B**, bottom) models demonstrate that  
 192 participants learned option values that recapitulate the correct ranking (i.e., from most negative to most  
 193 positive value) and are appropriately adaptive to the changes in outcome magnitudes beginning in Phase  
 194 3. Of note, VPRL produced more accurate estimates of the true state-action values of aversive options

195 (i.e., associated with monetary losses) over time compared to TDRL (**Figure S1**; two-way ANOVA (time,  
 196 model):  $F(\text{time}, 149) = 8.7$ ,  $p = 3.1e-83$ ;  $F(\text{model}, 1) = 66.2$ ,  $p = 2.4e-15$ ), whereas rewarding options are  
 197 learned with equivalent accuracy ( $F(\text{time}, 149) = 2.38$ ,  $p = 1.9e-13$ ;  $F(\text{model}, 1) = 0.02$ ,  $p = 0.66$ );  
 198 differences between learned option values for TDRL and VPRL models were specific to the most  
 199 negatively valued options.  
 200



**Figure 2 – TDRL and VPRL computational modeling results.** VPRL model best explains choice behavior on PRP task and leads to asymmetric learning. Distributions of (A) TDRL and (B) VPRL model parameter values across all participants. Horizontal bars mark the median of each distribution, and vertical bars indicate 95% credible interval (CI) of individual-level distribution. Dots indicate individual participant parameter values (mean of posterior parameter distributions); within-subject VPRL model parameters values are linked by grey lines. For both (A) and (B), bottom panels show time series of learned expected state-action values (Q-values) across participants for each option on the PRP task as predicted by the (A) TDRL and (B) VPRL models and shown relative to each option’s true expected value (grey dashed lines). Shaded ribbons indicate +/- one standard error of the mean (SEM), and different hues of shaded ribbons indicate different outcome probability groups. (C) Group- and individual-level differences in VPRL model learning rates (light blue distribution and histogram) and temporal discount parameters (dark blue distribution and histogram). Vertical dashed line indicates equivalence between parameter values; horizontal light and dark blue bars indicate 95% CI for the group-level distribution.

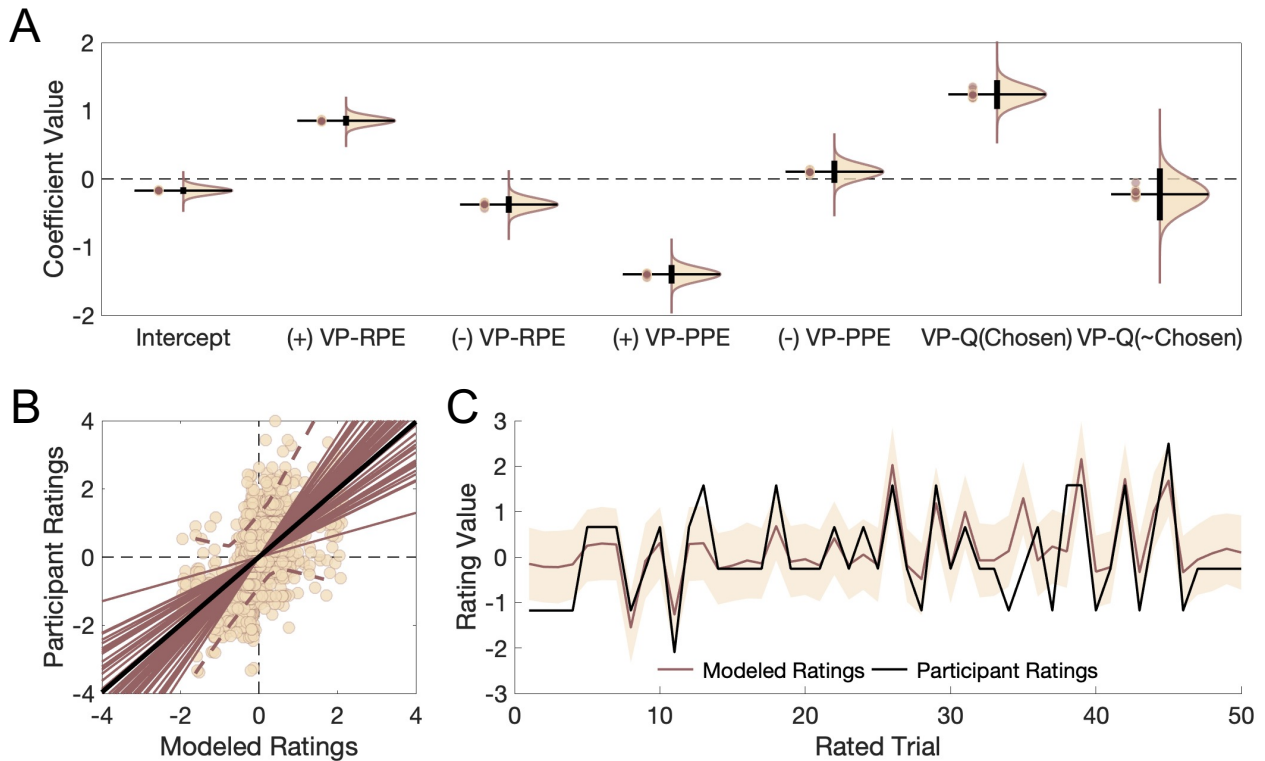
201  
 202  
 203  
 204  
 205  
 206  
 207  
 208  
 209  
 210  
 211  
 212  
 213  
 214  
 215  
 216  
 217

### VPRL prediction errors over rewards and punishments differentially affect subjective feelings

Given the evidence that VPRL best explains participant behavior on the PRP task, we sought to characterize how VPRL model-derived reward prediction errors (VP-RPE), punishment prediction errors (VP-PPE), and expected action values influence moment-to-moment changes in participants’ self-reported subjective feelings about experienced outcomes. We fit a cross-validated (leave-one-participant-out) Bayesian linear regression model to predict subjective feeling reports (**Figure 1**) using as predictor variables the learned state-action values (VP-Q-value) of each option presented on a rated trial and positive and negative VP-RPEs (from the VPRL Positive system) and VP-PPEs (from the VPRL Negative system) in response to the outcome of each rated trial (**Figure 3A**). Positive and negative VP-RPEs contribute positively and negatively to participants’ self-reported subjective feelings, respectively (positive VP-RPE: median = 0.92 [0.70 1.1]; negative VP-RPE: median = -0.33 [-0.64 -0.01]). Conversely, positive and negative VP-PPEs contribute negatively and positively to subjective feelings, respectively (positive VP-PPE: median = -1.4 [-1.8 -1.1]; negative punishment VP-PPE: median = 0.13 [0.02 0.21]).

218 The expected value of the chosen option on rated trials contributes positively to subjective ratings  
 219 (Expected Value (VP-EV) of Chosen: median = 1.4 [1.1 1.6]), whereas the expected value of the  
 220 unchosen option on rated trials shows no effect (VP-EV of Unchosen: median = -0.05 [-0.18 0.09]).  
 221

222 To assess the posterior predictive performance of the cross-validated Bayesian regression model  
 223 on held-out participant ratings, we computed r-squared values and Pearson correlation coefficients (rho  
 224 value) between the held-out participant's ratings and the model-predicted ratings (**Figure 3B**). This  
 225 analysis revealed a median within-participant correlation measure of 0.65 (SD = 0.16; median p-value =  
 226 4.3e-7) and r-squared value of 0.42 (SD = 0.18), indicating that the cross-validated regression model  
 227 generalizes moderately well to out-of-sample participant data (**Figure 3C**).  
 228



**Figure 3 – Dynamic changes in self-reported subjective feelings predicted by VPRL learning signals.** Cross-validated Bayesian regression analysis reveals influence of VPRL learning signals on ratings of subjective feelings about experienced outcomes. **(A)** Distribution of regression coefficient weights on positive (+) and negative (-) VP-RPEs and VP-PPEs and learned state-action values (VP-Q-values) of the chosen and unchosen options (VP-Q(Chosen) and VP-Q(~Chosen), respectively) on each trial. Horizontal bars indicate the median of each distribution; vertical bars indicate 95% CI of distribution; dots indicate mean individual parameter values. **(B)** Scatter plot demonstrating the linear relationship between model-derived and held-out participant ratings. Dark brown lines indicate within-participant linear correlations; black line indicates median linear relationship across all participant ratings; dashed brown lines represent 95% CI around individual-level linear correlation values. **(C)** Representative participant time series ( $\rho=0.77$ ,  $p = 7.9e-11$ ;  $r\text{-squared}=0.59$ ) of normalized subjective ratings (black line) and the cross-validated, model-predicted subjective ratings. Dark brown line represents the mean model-predicted ratings, and the shaded region represents  $\pm 1$  standard deviation around mean model predictions.

229  
 230  
 231 Further analyses comparing the coefficient values of VP-RPEs and VP-PPEs revealed that the  
 232 magnitude (absolute value) of positive VP-RPE coefficients is generally larger than the magnitude of  
 233 negative VP-RPE coefficients (median difference = 0.15 [-0.08 0.37]). Similarly, we found that positive  
 234 VP-PPEs have a consistently larger contribution to subjective ratings than negative VP-PPEs (median



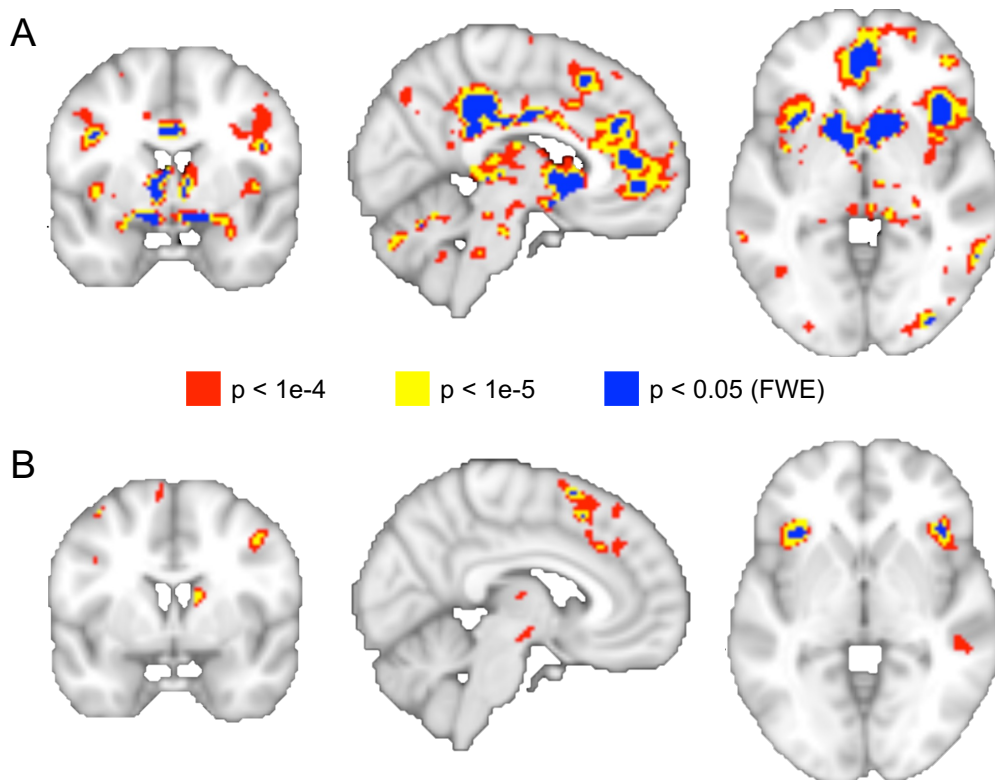
235 difference = 0.89 [0.63 1.1]). Lastly, negative VP-PPEs consistently had a more positive contribution to  
236 subjective ratings than negative VP-RPEs (median = 0.27 [-0.003 0.55]). As a whole, comparing the  
237 effects of positive and negative VP-RPEs and VP-PPEs on subjective ratings revealed a consistent  
238 ordering of the contributions of VPRL prediction errors to subjective feelings, with positive VP-RPE >  
239 negative VP-PPE (median difference = 0.98 [0.74 1.2]), negative VP-PPE > negative VP-RPE, and  
240 negative VP-RPE > positive VP-PPE (median difference = 0.62 [0.40 0.84]).

241  
242

### 243 ***Striatal-insular-prefrontal network activity tracks ‘reward’ and ‘punishment’ prediction errors***

244

245 Based on related prior work, we hypothesized that VPRL reward prediction errors (VP-RPEs),  
246 punishment prediction errors (VP-PPEs), and the respective positive and negative system expected  
247 values would be tracked by regions of dorsal and ventral striatum, cingulate cortex, and insula (O’Doherty  
248 et al., 2003; McClure et al., 2003; Pessiglione et al., 2006; Palminteri and Pessiglione, 2017; Seymour et  
249 al., 2004, 2005, 2007a, 2012; Delgado et al., 2008, 2009; Garrison et al., 2013). Using a model-based  
250 approach, we tested for regional activity that correlated with VP-RPEs and VP-PPEs by computing  
251 contrasts for positive effects of VP-RPEs or VP-PPEs (**Figure 4**). Regions that show hemodynamic  
252 signatures of VP-RPEs include the anterior cingulate cortex (ACC), anterior insula, and ventral striatum  
253 (**Figure 4A**); regions that correlate with VP-PPEs included the ACC, anterior insula, and dorsal striatum  
254 (**Figure 4B**). Additionally, we found that regional activity in the ventromedial prefrontal cortex (vmPFC)  
255 tracked VPRL-derived learned action values (VP-Q-values) of both the chosen and unchosen option on  
256 each trial (**Figure S2**).



**Figure 4 – Meso-cortico-limbic regional activity represents the set of VPRL prediction error signals.**

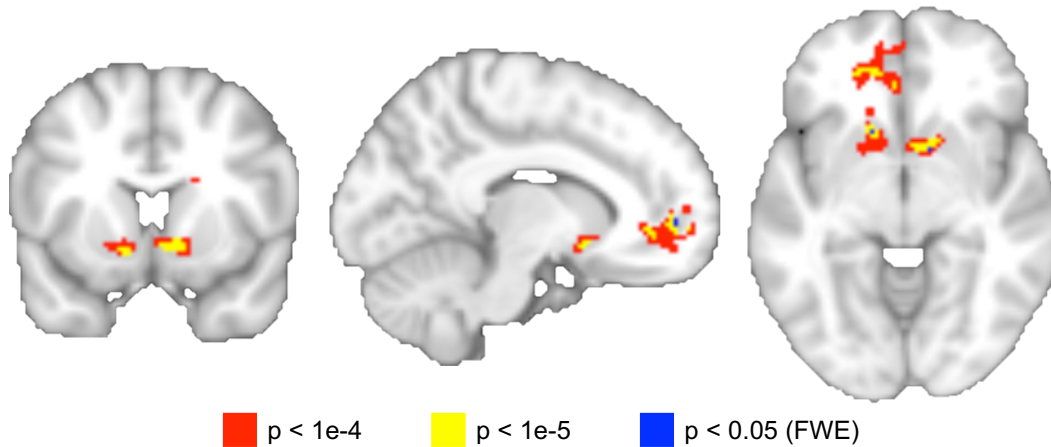
Whole-brain model-based analysis of VPRL learning signals indicate that regions of human striatum, insula cortex, and anterior cingulate cortex parametrically encode (A) VP-RPEs or (B) VP-PPEs. Colored voxels and associated p-values indicate statistical thresholding used for primary analyses. All panels are sliced at MNI coordinates  $x=6, y=2, z=-2$ . FWE = family-wise error.

257

258

259 **Ventral striatum and ventromedial prefrontal cortex track participants' subjective experience**

260  
261 Prior reports demonstrate that subjective feelings associated with prediction error and expected  
262 value signals are tracked by medial prefrontal cortex (Xiang et al., 2013) and ventral striatum (Rutledge  
263 et al., 2014). We hypothesized that the neural responses to signals derived from the VPRL model would  
264 drive brain responses associated with subjective feelings on every trial. Thus, we used the fitted model  
265 coefficients from the cross-validated subjective rating regression analysis to impute the subjective feeling  
266 on each trial conditioned on the trial's prediction errors and expected value signals. We found that  
267 regional activity in ventral striatum and ventromedial prefrontal cortex (vmPFC) parametrically tracked  
268 the imputed subjective rating on each trial throughout the PRP task (**Figure 5**).



**Figure 5 – Medial prefrontal and ventral striatal activity correlates of trial-to-trial subjective human feelings.** Whole-brain model-based analysis of self-reported subjective feelings as predicted by VPRL expected values and prediction errors indicate that medial prefrontal cortex and ventral striatum parametrically track the imputed subjective feeling on each trial. Colored voxels and associated p-value indicates statistical thresholding used for primary analyses;  $x=-10$ ,  $y=-10$ ,  $z=6$ . FWE = family-wise error.

269  
270  
271  
272  
273

274 **DISCUSSION**

275  
276 Here, we investigated how the human brain learns from independent appetitive and aversive  
277 experiences to adapt choice behaviors and how this dynamic process impacts subjective experience. We  
278 hypothesized VPRL (Kishida and Sands, 2021) as a framework for studying how neural systems may  
279 process rewards statistically independently from punishments. We demonstrate that VPRL consistently  
280 explains human choice behavior better on a probabilistic reward and punishment task compared to  
281 traditional TDRL. Furthermore, we show that VPRL-derived expected action values and prediction errors  
282 predict participants' self-reported ratings of subjective feelings to rewards and punishments trial-to-trial.  
283 Moreover, we demonstrated that these VPRL-derived learning signals are parametrically tracked by  
284 BOLD activity in dorsal and ventral striatum, cingulate cortex, anterior insula, and prefrontal cortex brain  
285 regions, and that BOLD signals in the ventral striatum and ventromedial prefrontal cortex track the  
286 expected rating of participants' subjective experience on each trial. Together, our results provide insight  
287 into (1) the type of learning mechanisms in humans responsible for ascribing valence information to  
288 stimuli and actions based on experience, and (2) how distributed neural activity implementing such  
289 mechanisms may support the composition of subjective phenomenal experience.

290  
291 Central to the VPRL hypothesis is the idea that there are parallel neural systems that process  
292 positive and negative experiences separately using TD learning before respective learning signals or  
293 valent value representations are available for further processing (Montague et al., 2016; Kishida and  
294 Sands, 2021). In this way, VPRL provides a novel computational framework for investigating a variety of  
295 neurophysiological, behavioral, and psychological phenomena. From an evolutionary perspective, early  
296 vertebrates likely developed neural circuits for detecting and escaping threatening (i.e., aversive) stimuli  
297 alongside, but separate from, (putatively) dopaminergic neural circuits for initial forms of associative  
298 reward learning (Cisek, 2021). This evolutionary theory is consistent with the idea that mammalian choice  
299 behavior may be driven by the activities of – and interaction between – separate positive and negative  
300 valence-processing systems, an idea with a venerable history in psychological theories of emotions  
301 (Cacioppo et al., 1999; Folkman and Moskowitz, 2000; Larsen et al., 2009; Larsen and McGraw, 2011)  
302 and motivated behaviors (Konorski, 1967; Dickinson and Dearing, 1979; Seymour et al., 2007b; Boureau  
303 and Dayan, 2013). Here, the VPRL framework can be viewed as an explicit generative account for the  
304 wide repertoire of 'approach-avoid' motivated behaviors (Dickinson and Dearing, 1979) and valence-  
305 specific affective responses (Cacioppo et al., 1999; Folkman and Moskowitz, 2000), while also providing  
306 a theoretical framework for considering the mechanisms of interaction between opponent systems that  
307 may lead to 'freezing', non-action responses (Boureau and Dayan, 2013), or the subjective phenomenal  
308 experience of 'mixed' or 'conflicting' emotions (Larsen and McGraw, 2011).

309  
310 In line with these evolutionary and psychological theories, we hypothesize a VPRL model that  
311 accounts for the premise that costs and benefits are always intertwined for biological creatures  
312 constrained by metabolic, survival, and reproductive goals and demands (Montague and King-Casas,  
313 2007; Botvinik et al., 2015). Importantly, VPRL specifies a different perspective on how valence, within  
314 the context of RL, is processed: aversive stimuli that are immediately (or predicted to be) costly are  
315 learned directly and independently from potential rewarding stimuli. This is distinct from the more  
316 common representations that requires aversive stimuli to be compared to 'expectations' and requires  
317 prediction error encoding according to the valence of the TD-RPE (i.e., differential learning from positive  
318 or negative RPEs).

319  
320 Our present results using a simple probabilistic reward and punishment learning task indicate  
321 that independently processing rewards and punishments via VPRL reveals an increased sensitivity to  
322 immediate punishments compared to rewards and increased temporal discounting of future punishments  
323 compared to future rewards (**Figures 2,3**), which is consistent with prior behavioral observations  
324 (Kahneman and Tversky, 1979; Tom et al., 2007). This differential learning from gains versus losses  
325 within the VPRL framework reveals that participants learn expected values of reward at a similar rate to  
326 that expected via traditional unidimensional TDRL, though they learn expected values of losses at a much  
327 faster rate and with improved accuracy and precision (**Figure S1**). This observation suggests that VPRL

328 signals may also independently and asymmetrically influence subjective feelings. Our results suggest  
329 that omitted or 'smaller-than-expected' rewards (i.e., negative VP-RPEs) do not contribute to what it 'feels  
330 like' in the same manner as 'larger-than-expected' punishments (positive VP-PPEs), nor do 'smaller-than-  
331 expected' punishments (negative VP-PPE) contribute to what it 'feels like' in the same manner as 'larger-  
332 than-expected' rewards (positive VP-RPEs). Such distinctions cannot be parsed within a unidimensional  
333 representation of valence as in the traditional TDRL framework. The presence of positive VP-RPEs had  
334 a significantly greater positive influence on subjective ratings than negative VP-PPEs; the opposite was  
335 also true: negative VP-RPEs had a consistently smaller negative influence on subjective ratings than  
336 positive VP-PPEs. Such relationships might reflect a relative scaling principle for positive and negative  
337 prediction errors over rewards and punishments in generating momentary affective subjective feelings,  
338 an effect that may be dependent on ventral striatal and ventromedial prefrontal neural activity (**Figure 5**).  
339 Regardless of the mechanisms to be discovered, our results demonstrate that VPRL is a valid  
340 neurocomputational framework for empirically investigating how complex interactions of reward and  
341 punishment may lead to self-reports about subjective phenomenal experience in humans.  
342

343 Numerous investigations into the neural basis of prediction error signaling in humans, using a  
344 variety of computational models and experimental designs, implicate a distributed network of brain  
345 regions in tracking prediction errors (Garrison et al., 2013). Our model-based event-related fMRI results  
346 indicate that VPRL prediction errors over rewards and punishments are represented by partially  
347 overlapping regional activation patterns and along a ventral-dorsal axis within the striatum (**Figures 4,5**),  
348 which is consistent with prior work (Seymour et al., 2007a). We hypothesize that striatal, insular,  
349 cingulate, and prefrontal cortex functional interactions – driven by an underlying neural architecture that  
350 broadcasts VPRL learning signals throughout the brain – can be viewed collectively as part of a dynamic  
351 affective core that regulates behavioral control and is a core component underlying subjective  
352 phenomenal experience (Kishida and Sands, 2021). Indeed, ascending neuromodulatory systems that  
353 project throughout both subcortical and cortical brain regions are integral to coordinating systems-level  
354 functional interactions towards accomplishing or switching between cognitive or behavioral tasks (Shine  
355 et al., 2019, 2021). Along these lines, future work may address how dynamic patterns of activity within  
356 the distributed subcortical-cortical network identified in our VPRL model-based analysis forms  
357 representations of state-action-outcome associations and how they co-evolve with representations of  
358 affective subjective experiences. In this regard, our results outline a potential role for ventral striatum and  
359 ventromedial prefrontal cortex interactions in mediating experience-dependent changes in brain activity  
360 associated with dynamic changes in subjective phenomenal experience (Xiang et al., 2013; Rutledge et  
361 al, 2014; Eldar et al., 2016; Tom et al., 2007; Chang et al., 2021), consistent with the dynamic affective  
362 core hypothesis (Kishida and Sands, 2021).  
363

364 Non-invasive brain activity measurements like fMRI are unable to provide information on the  
365 neurochemical basis of VPRL-reward prediction errors or VPRL-punishment prediction errors, though  
366 recent advances provide an opportunity for testing competing hypotheses (Kishida et al., 2016; Moran et  
367 al., 2018; Bang et al., 2020). Neurobiologically, an independent aversive system involved in valence  
368 processing may be implemented by a variety of possible neural substrates, such as a distinct population  
369 of dopamine neurons tuned for aversive stimuli (Matsumoto and Hikosaka, 2009; Lammel et al., 2014;  
370 Kishida et al., 2016; Kishida and Sands, 2021) or the serotonin neurotransmitter system (Daw et al.,  
371 2002; Boureau and Dayan, 2013; Montague et al., 2016; Moran et al., 2018; Kishida and Sands, 2021).  
372 For instance, direct electrochemical recordings of dopamine and serotonin microfluctuations in human  
373 striatum during a sequential investment task (Kishida et al., 2016; Moran et al., 2018) are consistent with  
374 the notion that these neurotransmitter systems can act as positive and negative valence-processing  
375 systems, respectively (Montague et al., 2016; Kishida and Sands, 2021). Further, dissociable effects of  
376 rewards and punishments on reversal learning have been linked to dopamine and serotonin transporter  
377 polymorphisms, respectively (den Ouden et al., 2013). Distributional RL (Dabney et al., 2020) has  
378 recently been demonstrated as a mechanism for representing a wide dynamic range of reward  
379 magnitude; still, it remains unclear how dopamine neurons come to achieve varying value prediction 'set  
380 points' as well as whether and how they might encode a distribution over aversive experiences.  
381 Alternatively, VPRL-like hypotheses for future investigation might address the potential distributional

382 coding of rewards and punishments in candidate neuromodulatory systems (e.g., dopamine and  
383 serotonin; Montague et al., 2016; Kishida and Sands, 2021; Moran et al., 2018; Bang et al., 2020).

384

385 Human behavior and subjective self-reports about associated phenomenal experiences, good  
386 and bad, are multidimensional. Prior work investigating computational models and associated BOLD  
387 imaging measurements of brain activity associated with subjective experience, mood, and subjective  
388 well-being utilized traditional unidimensional reinforcement learning models as a framework (Delgado et  
389 al., 2006; Rutledge et al., 2014; Eldar et al., 2016) and inspired the present work. Here, however, we  
390 demonstrate that a unidimensional reward prediction error is not enough (in contrast to arguments  
391 presented in Silver et al., 2020; Vamplew et al., 2021) to fully account for the dynamics of human choice  
392 behavior and associated subjective experiences and can even be detrimental when reward-associated  
393 actions also incur substantial physical costs or other negative externalities that cannot be disentangled  
394 with traditional TDRL (Elfwing and Seymour, 2017). Instead, our results using VPRL suggest that (at  
395 least) two valence dimensions are necessary, but this is almost certainly far from a complete depiction of  
396 the generative signals involved in experiences and behaviors associated with ‘what it is like’ to be (Nagel,  
397 1974). Our results are consistent with a need to account for appetitive and aversive input in parallel,  
398 though independently, such that the integration of these signals can be performed downstream of the  
399 systems that generate the error signals. As but one possible approach, VPRL maintains the  
400 computational advantages of TDRL, but also better accounts for information that biological agents must  
401 track (e.g., costly punishments or losses) that are often independent from co-occurring appetitive stimuli.  
402 We have taken an initial step to test VPRL as a hypothetical framework for investigating basic questions  
403 about how humans adapt their choice behavior and how associated signals may account for subjective  
404 phenomenal experiences. Our findings imply that new insights may be gained should VPRL (or other  
405 valence-partitioning models) be applied to computational psychiatric problems (Montague et al., 2012;  
406 Huys et al., 2016; Redish and Gordon, 2016; Brown et al., 2021) where subjective suffering and  
407 fundamental changes in adaptive behavior characterize severe challenges to mental health.

408

409 **ACKNOWLEDGEMENTS**

410 The authors would like to thank Read Montague and Terry Lohrenz for comments on an earlier draft of  
411 the manuscript. This work was supported by NIH-NIMH R01MH121099 (KTK), NIH-NIDA R01DA048096  
412 (KTK), NIH-NIMH R01MH124115 (KTK), NIH-NIDA P50DA006634 (KTK), NIH 5KL2TR001420 (KTK),  
413 NIH-NIDA F31DA053174 (LPS), and NIH-NIDA T32DA041349 (LPS). We would like to acknowledge the  
414 Translational Imaging Program (TIP) of the Wake Forest CTSI, which is supported by the National Center  
415 for Advancing Translational Sciences (NCATS), National Institutes of Health, through Grant Award  
416 Number UL1TR001420.

417

418

419 **AUTHOR CONTRIBUTIONS**

420 L.P.S. – Collected data; designed and performed data analysis; interpreted results, wrote, edited  
421 manuscript drafts, and approved final manuscript.

422 A.J. – Coded behavioral tasks; collected data; analyzed data; edited and approved final manuscript

423 R.E.J. – Collected data; analyzed data; edited and approved final manuscript

424 J.D.T. – Analyzed data; edited and approved final manuscript

425 K.T.K. – Conceived the study; designed experiments; supervised and guided data collection and analysis;  
426 interpreted results, wrote, edited manuscript drafts, and approved final manuscript.

427

428

429 **COMPETING INTERESTS**

430 The authors declare no competing interests.

431

432

433 **DATA AND CODE AVAILABILITY**

434 Anonymized individual-level participant behavioral task data and MRI data used in this study may be  
435 made available upon submission of a formal project outline from any qualified investigator to the  
436 corresponding author and subsequent approval by the corresponding author in line with data protection  
437 regulations of Wake Forest University School of Medicine Institutional Review Board (IRB). Custom-  
438 written analysis scripts for generating the behavioral and imaging results of this manuscript are  
439 maintained in a private github repository (*insert link upon acceptance*) that may be shared upon request  
440 from any qualified investigator to the corresponding author.

441

442

443 **METHODS**

444  
445 **Participants**

446  
447 A total of 47 participants (across two neuroimaging experiments (n=20; and n=27) were recruited  
448 from the local Winston-Salem community to complete the probabilistic reward and punishment (PRP)  
449 task. In the first fMRI cohort, participants (n=20; 16 female) were recruited from the community in  
450 Winston-Salem, NC. For the second fMRI cohort, participants (n=27; 19 female) were recruited as 'control  
451 participants' for an ongoing study. Recruitment of these participants in the second fMRI cohort was similar  
452 to the first fMRI cohort. However, consent to participate included repeated visits to be completed after an  
453 initial visit where the tasks completed include the PRP task as well as more extensive behavioral  
454 characterization after the PRP task was completed; the first visit in this ongoing study is similar to the  
455 visit completed by participants in the first fMRI cohort, except that *after* the completion of the PRP task  
456 with scanning participants underwent a more involved psychiatric evaluation process to properly control  
457 for the observational experimental group. Informed written consent was obtained from each participant,  
458 and the experiment was approved by the Institutional Review Board (IRB#'s: IRB00042265,  
459 IRB00054337, and IRB00056131) of Wake Forest University Health Sciences (WFUHS). All experiments  
460 were conducted at WFUHS.

461  
462 Three participants' subjective rating data were not used in the regression modeling analysis due  
463 to limited variability in the responses on the subjective rating assessment (i.e., choosing the same rating  
464 across 90% of rated trials). This results in n=44 participants for the combined fMRI cohort. From our  
465 leave-one-participant-out cross-validation approach, we computed Pearson's correlation coefficient ( $\rho$ )  
466 and r-squared values for the cross-validated model-predicted ratings (defined as the mean of samples of  
467 the posterior predictive density for the held-out participant's ratings) and actual held-out participant  
468 ratings. This procedure was iterated across participants, such that each individual acted as the held-out  
469 participant once. We used the fitted subjective rating model coefficients for each participant (i.e., the  
470 mean model coefficients for the cross-validated model iteration when the participant was held out) to  
471 impute a subjective rating for all trials for that participant, which we incorporate into the participant's first-  
472 level GLM in our model-based fMRI analysis.

473  
474 **Probabilistic Reward and Punishment (PRP) task experimental procedure**

475  
476 The PRP task (**Figure 1c**) is a 150-trial, two-choice monetary reward and punishment learning  
477 task, where chosen options are reinforced probabilistically with either monetary gains (or no gain) or  
478 monetary losses (or no loss). Six options (represented by fractal images) comprise the set of possible  
479 actions, with each option assigned to one of three outcome probabilities (25%, 50%, and 75%) and one  
480 of two outcome valences (monetary gain or loss); the assignment of options to outcome probabilities and  
481 valences is randomized across participants. The task proceeds through three phases. At the beginning  
482 of the experiment (Phase 1, trials 1-25), each trial starts with the presentation of two of the three possible  
483 'gain/no gain options, and participants are reinforced with either a monetary gain or nothing (\$1 or \$0)  
484 according to the chosen option's fixed probability. In Phase 2 (trials 26-75), the game introduces trials  
485 which present two of the three 'loss/no loss' options that result in either a monetary loss or nothing (-\$1  
486 or \$0) with fixed probabilities. There are 25 'gain/no gain' and 25 'loss/no loss' trials randomly ordered in  
487 Phase 2. In Phase 3 (trials 76-150), two options are presented randomly such that any trial may consist  
488 of two 'gain/no gain options, two 'loss/no loss' options, or one 'gain /no gain and one 'loss/no loss' option.  
489 Moreover, in Phase 3 the outcome magnitudes of all options change: the 25%, 50%, or 75% 'gain' options  
490 now payout \$2.50, \$1.50, and \$0.50 respectively, and the 25%, 50%, or 75% 'loss' options now lose -  
491 \$1.25, -\$0.75, and -\$0.25, respectively (dashed lines in **Fig. 1A**, bottom).

492  
493 A participant is presented with two options at the beginning of each trial, and they select an option  
494 at their own pace. The unchosen option disappears at the same time the chosen option is highlighted,  
495 and this screen lasts for three seconds. The outcome is then displayed for one second followed by a  
496 blank screen that lasts for a random time interval (defined by a Poisson distribution with  $\lambda = 3$  seconds)

497 before the next trial begins. After each trial, with probability 0.33, the blank screen following outcome  
 498 presentation is followed by a subjective feeling rating screen with the text “How do you feel about the last  
 499 outcome?”. Participants are asked to rate with a visual-digital scale (**Figure 1**) their feelings about the  
 500 last outcome, after which the blank screen reappears for another random interval before a new trial  
 501 begins.

502  
 503 **Temporal Difference Reinforcement Learning (Q-learning) model**

504  
 505 In the standard ‘unidimensional’ TDRL model (Sutton 1988; Watkins and Dayan, 1992; Sutton  
 506 and Barto, 1998), the expected value of a state-action pair  $Q(s_i, a)$ , where  $i$  indexes discrete time points  
 507 in a trial, is updated following selection of action  $a$  in state  $s_i$  according to the update rule:

508  

$$Q(s_i, a_i) \leftarrow Q(s_i, a_i) + \alpha \delta_i \quad \text{eq. 1}$$

509  
 510 where  $0 < \alpha < 1$  is a learning rate parameter that determines the weight prediction errors have on  
 511 updating expected values, and  $\delta_i$  is the TD reward prediction error term:

512  

$$\delta_i = [outcome_i + \gamma \max_a Q(s_{i+1}, \tilde{a})] - Q(s_i, a_i) \quad \text{eq. 2}$$

513  
 514 where  $outcome_i$  is the outcome (positive or negative) experienced in state  $s_i$  after taking action  $a_i$ ,  $0 <$   
 515  $\gamma < 1$  is a temporal discount parameter that discounts outcomes expected in the future relative to  
 516 immediate outcomes (i.e., a temporal discounting parameter), and  $\max_a Q(s_{i+1}, \tilde{a})$  is the maximum  
 517 expected value over all actions  $\tilde{a}$  afforded in the next state  $s_{i+1}$ . We defined the trials of the PRP task as  
 518 consisting of the set of  $i = \{1, 2, 3, 4\}$  event time points (1: options presented; 2: action taken; 3: outcome  
 519 presented; 4: (terminal) transition screen). We modeled participant choices ( $choice_t$ ) on each trial  $t$  of the  
 520 PRP task with a softmax choice policy (i.e., categorical logit choice model) that assigns probability to  
 521 choosing each of the two options presented on a trial according to the learned Q-values of the two options  
 522 present. For example, for a trial that presents option 2 and option 5, the corresponding Q-values  
 523  $Q(s_1, opt_2)$  and  $Q(s_1, opt_5)$  are used to compute the probability of selecting each option (e.g., option 2):

524  
 525  

$$P(choice_t = opt_2 \mid Q(s_1, opt_2), Q(s_1, opt_5)) = \frac{e^{Q(s_1, opt_2)/\tau}}{e^{Q(s_1, opt_2)/\tau} + e^{Q(s_1, opt_5)/\tau}} \quad \text{eq. 3}$$

526  
 527 where  $0 < \tau < 20$  is a choice temperature parameter that determines the softmax function slope and  
 528 parameterizes an exploration versus exploitation trade-off where higher temperature values lead to a  
 529 more uncertain distribution of choices and low temperature values allow choices to be attracted to higher  
 530 expected values.

531  
 532 **Valence-Partitioned Reinforcement Learning (VPRL) model**

533  
 534 For Valence-Partitioned RL (VPRL; Kishida and Sands, 2021), we extend the TDRL framework,  
 535 but separate ‘outcomes’ and how they are processed based on the valence of the input. VPRL treats  
 536 ‘Positive’ ( $P$ ) and ‘Negative’ ( $N$ ) input as though separate, parallel,  $P$ - and  $N$ -systems maintain a partition  
 537 between appetitive and aversive input throughout processing.  $P$ - and  $N$ -system Q-values are estimated  
 538 ( $Q^P, Q^N$ , respectively) independently, though each in a TDRL manner (see **eq. 4-7**). We model their  
 539 integration in the simplest manner (**eq. 8**) when value-based decisions must be made (Note: alternative  
 540 approaches for integrating these value estimates may be investigated in future work).

541  
 542  $P$ - and  $N$ -systems update via TD-prediction errors on every episode, but by valence specific rules  
 543 ( $P$ -system: **eq. 4** and  $N$ -system: **eq. 5**). The  $P$ -system only tracks rewarding (i.e., appetitive) outcomes  
 544 ( $outcome_i > 0$ , **eq. 4**) and the  $N$ -system only tracks punishing (i.e., aversive) outcomes ( $outcome_i < 0$ ,



545 **eq. 5**); both systems encode the opposite-valence outcomes and null outcomes similarly – as though no  
 546 outcome occurred.

547  
 548 For the  $P$ -system, the reward-oriented TD prediction error therefore is  
 549

$$\delta_i^P = \begin{cases} outcome_i + \gamma^P * \max_a Q^P(s_{i+1}, \tilde{a}) - Q^P(s_i, a_i) & \text{if } outcome_i > 0 \\ 0 + \gamma^P * \max_a Q^P(s_{i+1}, \tilde{a}) - Q^P(s_i, a_i) & \text{if } outcome_i \leq 0 \end{cases} \quad \text{eq. 4}$$

550  
 551 where  $0 < \gamma^P < 1$  is the  $P$ -system specific temporal discounting parameter (directly analogous to the  
 552 standard TDRL temporal discounting parameter).  
 553

554 The  $N$ -system similarly encodes a punishment-oriented TD prediction error term:  
 555

$$\delta_i^N = \begin{cases} |outcome_i| + \gamma^N * \max_a Q^N(s_{i+1}, \tilde{a}) - Q^N(s_i, a_i) & \text{if } outcome_i < 0 \\ 0 + \gamma^N * \max_a Q^N(s_{i+1}, \tilde{a}) - Q^N(s_i, a_i) & \text{if } outcome_i \geq 0 \end{cases} \quad \text{eq. 5}$$

556  
 557 where  $0 < \gamma^N < 1$  is the  $N$ -system temporal discounting parameter and  $|outcome_i|$  indicates the absolute  
 558 value of the outcome. The absolute value of the outcome is taken to be interpreted as though the system  
 559 only updates on aversive stimuli and does so based solely on the varying magnitudes.  
 560

561 The  $P$ - and  $N$ -systems prediction errors update expectations of future rewards or punishments of  
 562 an action, respectively, according to the standard TD-learning update rule but, again, for each system  
 563 independently:  
 564

$$Q^P(s_i, a_i) \leftarrow Q^P(s_i, a_i) + \alpha^P \delta_i^P \quad \text{eq. 6}$$

$$Q^N(s_i, a_i) \leftarrow Q^N(s_i, a_i) + \alpha^N \delta_i^N \quad \text{eq. 7}$$

565  
 566 where  $0 < \alpha^P < 1$  and  $0 < \alpha^N < 1$  are learning rates for the  $P$ - and  $N$ -systems,  $Q^P(s_i, a_i)$  is the expected  
 567 state-action value learned by the  $P$ -system, and  $Q^N(s_i, a_i)$  is the expected state-action value learned by  
 568 the  $N$ -system.  
 569

570 We compute a composite state-action value term for each action by contrasting the  $P$ - and  $N$ -  
 571 system Q-values,  
 572

$$Q(s_i, a_i) \leftarrow Q^P(s_i, a_i) - Q^N(s_i, a_i) \quad \text{eq. 8}$$

573  
 574 which is entered into the categorical logistic choice model (e.g., softmax policy, eq. 3) as for the TDRL  
 575 model above.  
 576

### 577 **TDRL and VPRL hierarchical model parameterization**

578  
 579 We specified a hierarchical structure to the TDRL and VPRL computational models to fit  
 580 participant choice behavior on the PRP task. Individual-level parameter values (e.g., learning rates) are  
 581 drawn from group-level distributions over each model parameter. This hierarchical modeling approach  
 582 accounts for dependencies between model parameters and biases individual-level parameter estimates  
 583 towards the group-level mean, thereby increasing reliability and certainty in parameter estimates,  
 584 improving model identifiability, and avoiding overfitting (Ahn et al., 2017). These hierarchical models  
 585 therefore cast individual participant parameter values as deviations from a group mean.  
 586

587 Formally, the joint posterior distribution  $P(\phi, \theta|y, M)$  over group-level parameters  $\phi$  and individual-  
 588 level parameters  $\theta$  for a given model  $M$  conditioned on the data from the cohort of participants  $y$  takes  
 589 the form  
 590

$$P(\mathbf{w}|y, M) = \frac{p(y|\mathbf{w}, M)p(\mathbf{w}|M)}{p(y|M)} \quad \text{eq. 9}$$

591  
 592 We simplify our notation to  $P(\mathbf{w}|y, M)$ , where  $\mathbf{w} = \{\phi, \theta\}$ ; here,  $P(y|\mathbf{w}, M)$  is the likelihood of choice data  
 593  $y$  conditioned on the model parameters and hyperparameters,  $P(y|M)$  is the marginal likelihood (model  
 594 evidence) of the data given a model, and  $P(\mathbf{w}|M)$  is the joint prior distribution over model parameters as  
 595 defined by the model, which can be decomposed into the product of the prior on individual-level model  
 596 parameters conditioned on the model hyper-parameters  $P(\theta|\phi, M)$  times the prior over hyper-parameters  
 597  $P(\phi|M)$ . We define the prior distributions for individual-level model parameters (e.g.,  $\theta_{TDRL} = \{\alpha, \tau, \gamma\}$  for  
 598  $M = \text{TDRL}$ ) and the hyper-priors of the means  $-\infty < \mu_{(\cdot)} < +\infty$  and standard deviations  $0 < \sigma_{(\cdot)} < +\infty$  of  
 599 the population-level parameter distributions (e.g.,  $\phi_{TDRL} = \{\mu_\alpha, \mu_\tau, \mu_\gamma, \sigma_\alpha, \sigma_\tau, \sigma_\gamma\}$ ) to be standard normal  
 600 distributions. We estimated all parameters in unconstrained space (e.g.,  $-\infty < \mu_\gamma < +\infty$ ) and use the  
 601 inverse Probit transform to map bounded parameters from unconstrained space to the unit interval  $[0,1]$   
 602 before scaling estimates by the parameter's upper bound:  
 603

$$\mu_\gamma \sim \text{Normal}(0,1) \quad \text{eq. 10}$$

$$\sigma_\gamma \sim \text{Normal}^+(0,1) \quad \text{eq. 11}$$

$$\boldsymbol{\tau}' \sim \text{Normal}(0,1) \quad \text{eq. 12}$$

$$\boldsymbol{\tau} = \text{Probit}^{-1}(\mu_\gamma + \sigma_\gamma * \boldsymbol{\tau}') * 20 \quad \text{eq. 13}$$

604  
 605 where bold terms indicate a vector of parameter values over participants. This non-centered  
 606 parameterization (Papaspiliopoulos et al., 2007) and inverse Probit transformation creates a uniform prior  
 607 distribution over individual-level model parameters between specified lower and upper bounds. Note that  
 608 for learning rate and temporal discount parameters, the scaling factor (upper bound) was set to 1,  
 609 whereas it was set to 20 for the choice temperature parameter. We used the Hamiltonian Monte Carlo  
 610 (HMC) sampling algorithm in the probabilistic programming language Stan via the R package *rstan* (v.  
 611 2.21.2; Carpenter et al., 2017) to estimate the joint posterior distribution over group- and individual-level  
 612 model parameters for the TDRL and VPRL models for both cohorts individually. For both models and  
 613 each cohort, we executed 12,000 total iterations (2,000 warm-up) on each of 3 chains for a total of 30,000  
 614 posterior samples per model parameter. We inspected chains for convergence by verifying sufficient  
 615 chain mixing according to the Gelman-Rubin statistic  $\hat{R}$ , which was approximately 1 for all parameters.  
 616

### 617 TDRL and VPRL model comparison

618  
 619 We compared the TDRL and VPRL models' fits to participant choice behavior on the PRP task  
 620 according to their model evidence (i.e., model marginal likelihood), which represents the probability or  
 621 'plausibility' of observing the actual PRP task data under each model (Mckay 2013). In Bayesian model  
 622 comparison, the model with the greatest posterior model probability  $p(M|y)$  is deemed the best  
 623 explanation for the data  $y$  and is computed by:  
 624

$$P(M|y) \propto P(y|M)P(M) \quad \text{eq. 14}$$

625  
 626 where  $P(y|M)$  is the model marginal likelihood ("model evidence") and  $P(M)$  is the model's prior  
 627 probability. The model evidence is defined as:  
 628

$$P(y|M) = \int P(y|\mathbf{w}, M)P(\mathbf{w}|M)d\mathbf{w} \quad \text{eq. 15}$$

629

630 where  $P(\mathbf{w}|M)$  is the prior probability of a model  $M$ 's parameters  $\mathbf{w}$  before observing any data and  
 631  $P(y|\mathbf{w},M)$  is the likelihood of data  $y$  given a model and its parameters. We adopt the approach of  
 632 approximating this integral using importance sampling (i.e., bridge sampling). Given that we only wanted  
 633 to compare the TDRL and VPRL models, the relative posterior model probability can be defined as:  
 634

$$\frac{P(TDRL|y)}{P(VPRL|y)} = \frac{P(TDRL) * P(y|TDRL)}{P(VPRL) * P(y|VPRL)} \quad \text{eq. 16}$$

635  
 636 where the ratio of posterior model probabilities  $\frac{P(TDRL|y)}{P(VPRL|y)}$  is referred to as the “posterior  
 637 odds” of TDRL relative to VPRL;  $P(TDRL)$  and  $P(VPRL)$  are the prior probabilities of the TDRL and VPRL  
 638 models, respectively; and the ratio of marginal likelihoods  $\frac{P(y|TDRL)}{P(y|VPRL)}$  is termed the “Bayes  
 639 factor”, which is a standard measure for Bayesian model comparison. Granting equal prior probabilities  
 640 over the set of candidate models, each model's evidence  $P(y|M)$  can be used to rank each model in the  
 641 set for comparison. The marginal likelihoods are computed as log-scaled. We estimated the log model  
 642 evidence for the TDRL and VPRL models for each cohort using an adaptive importance sampling routing  
 643 called bridge sampling as implemented in the R package *bridgesampling* (v. 1.1-2; Gronau et al., 2017).  
 644 Bridge sampling is an efficient and accurate approach to calculating normalizing constants like the  
 645 marginal likelihood of models even with hierarchical structure and for reinforcement learning models in  
 646 particular (Gronau et al., 2017). To further ensure stability in the bridge sampler's estimates of model  
 647 evidence, we performed 10 repetitions of the sampler and report the median and interquartile range of  
 648 the estimates of model evidence. The model with the maximum (i.e., less negative) model evidence is  
 649 preferred, and therefore a positive value for the difference between the log model evidences for TDRL  
 650 and VPRL (as reported in **Table 1**) favors TDRL, while a negative value favors VPRL.  
 651

652 In addition to the standard Bayesian model comparison using model marginal likelihoods, we  
 653 estimated each model's Bayesian leave-one-out (LOO) cross-validation predictive accuracy, defined as  
 654 a model's expected log predictive density (ELPD-LOO):  
 655

$$elpd_{LOO} = \sum_{i=1}^N \log(p(y_i|y_{-i})) \quad \text{eq. 17}$$

656  
 657 where the posterior predictive distribution  $p(y_i|y_{-i})$  for held-out data  $y_i$  given a set of training data  $y_{-i}$ , is  
 658

$$P(y_i|y_{-i}) = \int p(y_i|\mathbf{w})p(\mathbf{w}|y_{-i})d\mathbf{w} \quad \text{eq. 18}$$

659  
 660 The ELPD is an estimate of (i.e., approximation to) the cross-validated accuracy of the TDRL or VPRL  
 661 models in predicting new (i.e., held-out) participant data, given the posterior distribution over model  
 662 parameters fit to a training set of participant data (Vehtari et al., 2017). Again, we approximate this integral  
 663 via importance sampling of the joint posterior parameter distribution given the training data  $p(\mathbf{w}|y_{-i})$ .  
 664 Furthermore, the upper tail of the distribution of importance weights are smoothed by a Pareto distribution  
 665 (Pareto-smoothed importance sampling, PSIS) to improve the ELPD-LOO estimation. We calculated the  
 666 model ELPD in this way using the R package *loo* (v. 2.3.1; Vehtari et al., 2017).  
 667

## 668 **Subjective rating computational modeling and cross-validated Bayesian regression analysis**

669  
 670 We defined a Bayesian linear regression model of Positive and Negative valence system  
 671 prediction errors and estimated Q-values on participants' self-reported subjective feelings about their  
 672 most recent outcomes measured throughout the PRP task (query probability = .33 on each trial). We  
 673 express the subjective rating on a trial as normally distributed with a mean  $E(y_i|\beta, X)$  that is a linear

674 function of Positive and Negative system prediction errors and learned Q-values (predictor variable matrix  
 675  $X$ ):

$$E(y_i|\beta, X) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon_i \quad \text{eq. 19}$$

$$\varepsilon \sim \text{Normal}(0, \sigma^2) \quad \text{eq. 20}$$

677  
 678 where  $E(y_i)$  is the subjective rating ( $i = 1 \dots 50$  indexes the numbers of ratings) from a participant,  $\beta_k$  are  
 679 the  $k = 7$  linear model weights,  $x_k$  are rows of the predictor variable matrix  $X$  (corresponding to each trial  
 680 on which a subjective rating was sampled), and  $\varepsilon_i$  are the normally-distributed errors with variance  $\sigma^2$ .  
 681 We define  $\theta = \{\beta_0, \beta_1, \dots, \beta_k, \sigma\}$  as the vector of all model parameters. The Bayesian rendering of the  
 682 subjective rating linear regression model is therefore

$$p(\theta|y, X) \propto p(y|\theta, X)p(\theta) \quad \text{eq. 21}$$

684  
 685 where  $p(y|\theta, X)$  is the (normally-distributed) data likelihood function and  $p(\theta) = p(\beta)p(\sigma^2)$  are the  
 686 (weakly-informative) prior distributions over model parameters:

$$p(y|\theta, X) \sim \text{Normal}(X\beta, \sigma^2 I) \quad \text{eq. 22}$$

$$p(\beta) \sim \text{Normal}(0, 1) \quad \text{eq. 23}$$

$$p(\sigma^2) \sim \text{Exponential}(1) \quad \text{eq. 24}$$

688  
 689 where  $I$  is the  $n \times n$  ( $n =$  number of participants in training sample) identity matrix. The joint posterior  
 690 distribution over model parameters  $\theta$  conditioned on the subjective ratings and predictor variable matrix  
 691 factorizes into:

$$p(\theta|y, X) \propto p(\beta|\sigma^2, y, X)p(\sigma^2|y, X) \quad \text{eq. 25}$$

693  
 694 where the conditional posterior distribution  $p(\beta|\sigma^2, y, X)$  of linear model parameters  $\beta$  conditional on  $\sigma^2$   
 695 is the normal distribution

$$p(\beta|\sigma^2, y, X) \sim \text{Normal}(\hat{\beta}, V_B \sigma^2) \quad \text{eq. 26}$$

697  
 698 and, from the least-squares solution,

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad \text{eq. 27}$$

$$V_B = (X^T X)^{-1} \quad \text{eq. 28}$$

700  
 701 The marginal posterior distribution  $p(\sigma^2|y, X)$  is defined as

$$p(\sigma^2|y, X) \sim \text{Inverse} - \chi^2(n - k, s^2) \quad \text{eq. 29}$$

$$s^2 = \frac{1}{n - k} (y - X\hat{\beta})^T (y - X\hat{\beta}) \quad \text{eq. 30}$$

703  
 704 where  $n - k$  is the number of degrees of freedom (data points). We implemented this Bayesian regression  
 705 model using the R package *rstanarm* (v. 2.21.1; Gabry and Goodrich, 2017), which uses HMC via Stan  
 706 to efficiently sample the entire joint posterior distribution over model parameters  $p(\theta|y, X)$ . We adopted  
 707 a leave-one-participant-out cross-validation approach by fitting the subjective rating regression model to  
 708 all participants except for one person and drawing samples of  $(\beta, \sigma)$  from this fitted model's joint posterior  
 709 distribution to form a posterior predictive distribution  $p(\tilde{y}|y)$  for the held-out participant's ratings  $\tilde{y}$  as:

$$p(\tilde{y}|y) \sim \text{Normal}(\tilde{X}\beta, \sigma^2 I) \quad \text{eq. 31}$$

711

712 where  $\tilde{X}$  is the held-out participant's predictor matrix. For sampling both the linear model joint posterior  
713 distribution  $p(\theta|y, X)$  and the posterior predictive distribution  $p(\tilde{y}|y)$ , we drew 3,500 total samples (1,000  
714 warm-up) on each of 4 chains for a total of 10,000 samples for each parameter and verified sufficient  
715 mixing according to  $\hat{R}$  values, which were approximately 1 for all parameters.

716

### 717 **Functional MRI data acquisition, pre-processing, and model-based analysis**

718

719 The fMRI cohort (N=47) was recruited as part of two separate studies at WFUHS, with one cohort  
720 n=20 and the other n=27. For all neuroimaging participants, we acquired fMRI and structural MRI data  
721 using a Siemens MAGNETOM 3T Skyra whole-body scanner and a 32-channel head coil. High-resolution  
722 (0.5x0.5x1.0mm<sup>3</sup>) T1-weighted structural MRI scans were acquired using a magnetization-prepared rapid  
723 gradient echo (MPRAGE) sequence (TR = 1480msec ; TE = 2.66msec ; flip angle = 12 degrees; FoV =  
724 24.5cm ; 192 slices), and fMRI BOLD data were acquired by means of a multi-band (simultaneous  
725 multislice, SMS) echo-planar imaging (EPI) sequence (MB factor = 8; TR = 1000msec; TE = 30msec; flip  
726 angle = 52 degrees; FoV = 20.8cm; 72 interleaved sagittal slices; isotropic 2mm<sup>3</sup> voxels). All data were  
727 pre-processed and analyzed using FSL and SPM12. Each participant's fMRI data were aligned to a  
728 single-band reference image (SBref) and corrected for EPI (B0) distortions using a fieldmap estimated  
729 from reverse-phase encoded functional volumes (directions Right->Left and Left->Right) via FSL's *topup*  
730 tool (Andersson et al., 2003); co-registered to the high-resolution structural volume and warped to MNI  
731 template space (2mm<sup>3</sup> isotropic); spatially smoothed with a 5mm FWHM Gaussian kernel; high-pass  
732 filtered at 128 seconds (<0.008Hz); and normalized by the session grand-mean value.

733

734 For each participant, we constructed a first-level GLM to model BOLD signals during task  
735 performance. The following regressors were included in the GLM as events of interest and convolved  
736 with a canonical hemodynamic response function: (i) onset of 'option presentation', parametrically  
737 modulated by (a) expected value of the chosen option and (b) expected value of the unchosen option;  
738 (ii) onset of 'outcome presentation', parametrically modulated by the (a) 'outcome presentation'-episode-  
739 specific positive system prediction error, (b) 'outcome presentation'-episode-specific negative system  
740 prediction error, and (c) imputed rating of subjective feelings; and (iii) all other motor and visual stimuli.  
741 The parametric modulators at the time of outcome presentation were orthogonalized. Six head motion  
742 parameters were included as covariates of no interest. First-level GLM results for each participant were  
743 incorporated into a second-level random effects analysis at the group-level. At the group-level, all  
744 analyses were whole-brain and conducted at either a family-wise error (FWE)-corrected statistical  
745 threshold of  $p < 0.05$  or an uncorrected significance thresholds of  $p < 1e-4$  and  $p < 1e-5$ .

746

747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800

## REFERENCES

- Bang, D., Kishida, K.T., Lohrenz, T., White, J.P., Laxton, A.W., Tatter, S.B., Fleming, S.M., and Montague, P. R. (2020). Sub-second dopamine and serotonin signaling in human striatum during perceptual decision-making. *Neuron* 108, 999-1010.
- Bayer, H.M., and Glimcher, P.W. (2005). Midbrain dopamine neurons encode a quantitative reward prediction error signal. *Neuron* 47, 129-141.
- Bayer, H.M., Lau, B., and Glimcher, P.W. (2007). Statistics of midbrain dopamine neuron spike trains in the awake primate. *Journal of Neurophysiology* 98(3), 1428-1439.
- Botvinick, M., Weinstein, A., Solway, A., and Barto, A. (2015). Reinforcement learning, efficient coding, and the statistics of natural tasks. *Current Opinion in Behavioral Sciences* 5, 71-77.
- Brown, V.M., Zhu, L., Solway, A., Wang, J.M., McCurry, K.L., King-Casas, B., and Chu, P. (2021). Reinforcement learning disruptions in individuals with depression and sensitivity to symptom change following cognitive behavioral therapy. *JAMA Psychiatry* 78(10), 1113-1122.
- Cacioppo, J.T., Gardner, W.L., and Berntson, G.G. (1999). The affect system has parallel and integrative processing components: form follows function. *Journal of Personality and Social Psychology* 76(5), 839-855.
- Chalmers, D. (1995). Facing up to the problem of consciousness. *Journal of Consciousness Studies* 2(3), 200-219.
- Churchland, P.M. (1984). *Matter and Consciousness* (MIT Press).
- Churchland, P.M. (2014). Consciousness and the introspection of qualitative simples. In *Consciousness Inside and Out: Phenomenology, Neuroscience, and the Nature of Experience*, R. Brown, ed. (Springer), pp. 35-56.
- Churchland, P.S. (1996). The hornswoggle problem. *Journal of Consciousness Studies* 3(5-6), 402-408.
- Churchland, P. S., & Sejnowski, T. J. (1994). *The computational brain* (MIT Press).
- Chang, L.J., Jolly, E., Cheong, J.H., Rapuano, K.M., Greenstein, N., Chen, P.-H.A., and Manning, J.R. (2021). Endogenous variation in ventromedial prefrontal cortex state dynamics during naturalistic viewing reflects affective experience. *Science Advances* 7(17), 1-17.
- Cisek, P. (2021). Evolution of behavioral control from chordates to primates. *Philosophical Transactions of the Royal Society B* 377, 20200522.
- Dabney, W., Kurth-Nelse, Z., Uchida, N., Starkweather, C.K., Hassabis, D., Munos, R., and Botvinick, M. (2020). A distributional code for value in dopamine-based reinforcement learning. *Nature* 577(7792), 671-675.
- Daw, N. D., Kakade, S., & Dayan, P. (2002). Opponent interactions between serotonin and dopamine. *Neural Networks* 15(4), 603-616.
- Dayan, P., and Niv, Y. (2008). Reinforcement learning: The Good, The Bad and The Ugly. *Current Opinion in Neurobiology* 18(2), 185-196.
- Delgado, M.R., Labouliere, C.D., and Phelps, E.A. (2006). Fear of losing money? Aversive conditioning with secondary reinforcers. *Social, Cognitive, and Affective Neuroscience* 1(3), 250-259.
- Delgado, M.R., Li, J., Schiller, D., and Phelps, E.A. (2008). The role of the striatum in aversive learning and aversive prediction errors. *Philosophical Transactions of the Royal Society B: Biological Sciences* 363(1511), 3787-3800.
- Delgado, M.R., Jou, R.L., LeDoux, J.E., and Phelps, E.A. (2009). Avoidance of negative outcomes: tracking the mechanisms of avoidance learning in humans during fear conditioning. *Frontiers in Behavioral Neuroscience* 3(33), 1-9.
- den Ouden, H.E.M., Friston, K.J., Daw, N.D., McIntosh, A.R., and Stephan, K.E. (2013). Dissociable effects of dopamine and serotonin on reversal learning. *Neuron* 80, 1090-1100.
- Dickinson, A., and Dearing, M.F. (1979). Appetitive-aversive interactions and inhibitory processes. In *Mechanisms of Learning and Motivation*, A. Dickinson and R.A. Boakes, eds. (Psychology Press), pp. 203-231.
- Eldar, E., Rutledge, R.B., Dolan, R.J., & Niv, Y. (2016). Mood as Representation of Momentum. *Trends in Cognitive Sciences* 20(1), 15-24.
- Elfving, S., and Seymour, B. (2017). Parallel reward and punishment control in humans and robots: Safe reinforcement learning using the MaxPain algorithm. *Joint IEEE International Conference on Development and Learning and Epigenetic Robotics (ICDL-EpiRob)*, pp. 140-147.

801 Eshel, N., Bukwich, M., Rao, V., Hemmelder, V., Tian, J., and Uchida, N. (2015). Arithmetic and local  
802 circuitry underlying dopamine prediction errors. *Nature* 573(7568), 243-246.

803 Faherty, A., ed. (2016). *States of Mind: Experiences at the Edge of Consciousness: A Collection of*  
804 *Literature, Science, Philosophy and Art*. Wellcome Collection, part of The Wellcome Trust.

805 Folkman, S., and Moskowitz, J.T. (2000). Positive affect and the other side of coping. *American*  
806 *Psychologist* 55(6), 647-654.

807 Forbes, N., & Mahon, B. (2014). *Faraday, Maxwell, and the electromagnetic field: How two men*  
808 *revolutionized physics* (Prometheus Books).

809 Garrison, J., Erdeniz, B., and Done, J. (2013). Prediction error in reinforcement learning: a meta-analysis  
810 of neuroimaging studies. *Neuroscience & Biobehavioral Reviews* 37(7), 1297-1310.

811 Glimcher, P. (2011). Understanding dopamine and reinforcement learning: the dopamine reward  
812 prediction error hypothesis. *Proceedings of the National Academy of Sciences* 108, 15647-15654.

813 Hart, A.S., Rutledge, R.B., Glimcher, P.W., and Phillips, P.E.M. (2014). Phasic dopamine release in rat  
814 nucleus accumbens symmetrically encodes a reward prediction error. *Journal of Neuroscience*  
815 34(3), 698-704.

816 Huys, Q.J.M., Maia, T.V., and Frank, M.J. (2016). Computational psychiatry as a bridge from  
817 neuroscience to clinical applications. *Nature Neuroscience* 19, 404-413.

818 Kahneman, D., and Tversky, A. (1979). Prospect theory: an analysis of decision under risk. *Econometrica*  
819 47(2), 263-292.

820 Kishida, K.T., King-Casas, B., and Montague, P.R. (2010). Neuroeconomic Approaches to Mental  
821 Disorders. *Neuron* 67, 543-554. PMID: 20797532.

822 Kishida, K.T. (2012). A computational approach to “free will” constrained by the games we play. *Frontiers*  
823 *in Integrative Neuroscience* 6. PMID: 23060761

824 Kishida, K. T., Saez, I., Lohrenz, T., Witcher, M.R., Laxton, A.W., Tatter, S.B., White, J.P., Ellis, T.L.,  
825 Phillips, P. E.M., and Montague, P.R. (2016). Subsecond dopamine fluctuations in human  
826 striatum encode superposed error signals about actual and counterfactual reward. *Proceedings*  
827 *of the National Academy of Sciences of the United States of America* 113(1), 200–205.

828 Kishida, K.T., and Sands, L.P. (2021). A dynamic affective core to bind the contents, context, and value  
829 of conscious experience. In *Affect Dynamics*, C. Waugh and P. Kuppens, eds. (Springer), pp.  
830 293-328.

831 Konorski, J. (1967). *Integrative Activity of the Brain* (University of Chicago Press: Chicago).

832 Lammel, S., Lim, B.K., and Malenka, R.C. (2014). Reward and aversion in a heterogeneous midbrain  
833 dopamine system. *Neuropharmacology* 76, 351-359.

834 Larsen, J.T., Norris, C.J., McGraw, A.P., Hawkey, L.C., and Cacioppo, J.T. (2009). The evaluative space  
835 grid: A single-item measure of positivity and negativity. *Cognition and Emotion* 23(3), 453-480.

836 Larsen, J.T., and McGraw, A.P. (2011). Further evidence for mixed emotions. *Journal of Personality and*  
837 *Social Psychology* 100(6), 1095-1110.

838 Lefebvre, G., Lebreton, M., Florent, M., Bourgeois-Gironde, S., and Palminteri, S. (2017). Behavioral and  
839 neural characterization of optimistic reinforcement learning. *Nature Human Behavior* 1(4), 1-9.

840 Matsumoto, M., and Hikosaka, O. (2009). Two types of dopamine neurons distinctly convey positive and  
841 negative motivational signals. *Nature* 459(7248), 837-841.

842 McClure, S.M., Berns, G.S., and Montague, P.R. (2003). Temporal prediction errors in a passive learning  
843 task activate human striatum. *Neuron* 38(2), 339-346.

844 Montague, P.R., Dayan, P., and Sejnowski, T.J. (1996). A framework for mesencephalic dopamine  
845 systems based on predictive Hebbian learning. *Journal of Neuroscience* 16(5), 1936–1947.

846 Montague, P.R., and King-Casas, B. (2007). Efficient statistics, common currencies and the problem of  
847 reward harvesting. *Trends in Cognitive Sciences* 11(12), 514-519.

848 Montague, P.R., Dolan, R.J., Friston, K.J., and Dayan, P. (2012). Computational psychiatry. *Trends in*  
849 *Cognitive Sciences* 16(1), 72-80.

850 Montague, P.R., Kishida, K.T., Moran, R.J., and Lohrenz, T.M. (2016). An efficiency framework for  
851 valence processing systems inspired by soft cross-wiring. *Current Opinion in Behavioral Sciences*  
852 11, 121–129.

853 Moran, R.J., Kishida, K.T., Lohrenz, T., Saez, I., Laxton, A.W., Witcher, M.R., Tatter, S.B., Ellis, T.L.,  
854 Phillips, P.E., Dayan, P., and Montague, P.R. (2018). The Protective Action Encoding of Serotonin  
855 Transients in the Human Brain. *Neuropsychopharmacology* 43(6), 1425–1435.

856 Nagel, T. (1974). What is it like to be a bat? *Philosophical Review* 83(4), 435-450.

857 O’Doherty, J.P., Dayan, P., Friston, K., Critchley, H., and Dolan, R.J. (2003). Temporal difference models  
858 and reward-related learning in the human brain. *Neuron* 38(2), 329–337.

859 Palminteri, S., and Pessiglione, M. (2017). Opponent brain systems for reward and punishment learning:  
860 causal evidence from drug and lesion studies in humans. In *Decision Neuroscience: An*  
861 *Integrative Perspective*, J.C. Dreher and L. Tremblay, eds. (Academic Press, San Diego), pp.  
862 291-303.

863 Pessiglione, M., Seymour, B., Flandin, G., Dolan, R.J., and Frith, C.D. (2006). Dopamine-dependent  
864 prediction errors underpin reward-seeking behavior in humans. *Nature* 442(7106), 1042-1045.

865 Pessiglione, M., and Delgado, M.R. (2015). The good, the bad and the brain: Neural correlates of  
866 appetitive and aversive values underlying decision making. *Current Opinion in Behavioral*  
867 *Sciences* 5, 78–84.

868 Redish, D. (2004). Addiction as a computational process gone awry. *Science* 306, 1944-1947.

869 Redish, D., and Gordon, J. (2016). *Computational Psychiatry: New Perspectives on Mental Illness* (MIT  
870 Press).

871 Rutledge, R.B., Skandali, N., Dayan, P., and Dolan, R.J. (2014). A computational and neural model of  
872 momentary subjective well-being. *Proceedings of the National Academy of Sciences* 111(33),  
873 12252–12257.

874 Rutledge, R.B., Moutoussis, M., Smittenaar, P., Zeidman, P., Taylor, T., Hrynkiewicz, L., Lam, J.,  
875 Skandali, N., Siegel, J.Z., Ousdal, O.T., et al. (2017). Association of neural and emotional  
876 impacts of reward prediction errors with major depression. *JAMA Psychiatry* 74(8), 790-797.

877 Schultz, W., Dayan, P., and Montague, P.R. (1997). A neural substrate of prediction and reward. *Science*  
878 275(5306), 1593–1599.

879 Seymour, B., O’Doherty, J.P., Dayan, P., Koltzenburg, M., Jones, A.K., Dolan, R.J., Friston, K.J., and  
880 Frackowiak, R.S. (2004). Temporal difference models describe higher-order learning in humans.  
881 *Nature* 429(6992), 664-667.

882 Seymour, B., O’Doherty, J.P., Koltzenburg, M., Wiech, K., Frackowiak, R., Friston, K., and Dolan, R.  
883 (2005). Opponent appetitive-aversive neural processes underlie predictive learning of pain relief.  
884 *Nature Neuroscience* 8(9), 1234–1240.

885 Seymour, B., Daw, N., Dayan, P., Singer, T., and Dolan, R. (2007a). Differential Encoding of Losses and  
886 Gains in the Human Striatum. *Journal of Neuroscience* 27(18), 4826–4831.

887 Seymour, B., Singer, T., and Dolan, R. (2007b). The neurobiology of punishment. *Nature Reviews*  
888 *Neuroscience* 8, 300-311.

889 Seymour, B., Daw, N., Roiser, J.P., Dayan, P., and Dolan, R. (2012). Serotonin selectively modulates  
890 reward value in human decision-making. *Journal of Neuroscience* 32(17), 5833-5842.

891 Shine, J.M., Müller, E.J., Munn, B., Cabral, J., Moran, R.J., and Breakspear, M. (2021). Computational  
892 models link cellular mechanisms of neuromodulation to large-scale neural dynamics. *Nature*  
893 *Neuroscience* 24, 765-776.

894 Shine, J.M., Breakspear, M., Bell, P.T., Martens, K.A.E., Shine, R., Koyejo, O., Sporns, O., and Poldrack,  
895 R.A. (2019). Human cognition involves the dynamic integration of neural activity and  
896 neuromodulatory systems. *Nature Neuroscience* 22, 289-296.

897 Silver, D., Singh, S., Precup, D., and Sutton, R.S. (2021). Reward is enough. *Artificial Intelligence* 299,  
898 1-13.

899 Sutton, R.S. (1988). Learning to predict by the methods of temporal differences. *Machine Learning* 3(1),  
900 9-44.

901 Sutton, R.S., and Barto, A. (1998). *Reinforcement Learning: An Introduction* (MIT Press).

902 Taschereau-Dumouchel, V., Michel, M., Lau, H., Hofmann, S.G., and Ledoux, J.E. (2022). Putting the  
903 “mental” back in “mental disorders”: a perspective from research on fear and anxiety. *Molecular*  
904 *Psychiatry*, In Press.

905 Tom, S.M., Fox, C.R., Trepel, C., and Poldrack, R.A. (2007). The neural basis of loss aversion in decision-  
906 making under risk. *Science* 315, 515-518.



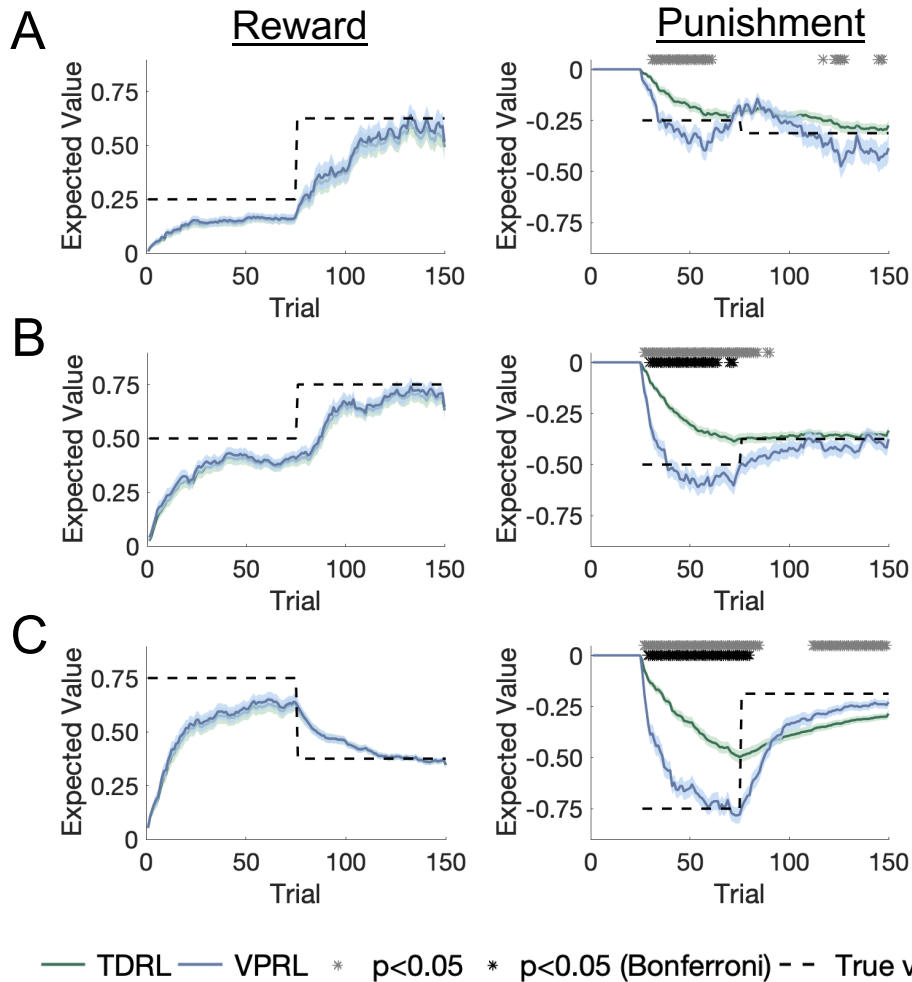
- 907 Vamplew, P., Smith, B.J., Kallstrom, J., Ramos, G., Radulescu, R., Roijers, D.M., Hayes, C.F., Heintz,  
908 F., Mannion, P., Libin, P.J.K., et al. (2021). Scalar reward is not enough: a response to Silver,  
909 Singh, Precup, and Sutton (2021). arXiv:2112.15422.
- 910 Watabe-Uchida, M., Eshel, N., and Uchida, N. (2017). Neural circuitry of reward prediction error. *Annual*  
911 *Review of Neuroscience* 40, 373-394.
- 912 Xiang, T., Lohrenz, T., and Montague, P.R. (2013). Computational substrates of norms and their  
913 violations during social exchange. *Journal of Neuroscience* 33(3), 1099-1108.
- 914 Zaghoul, K.A., Blanco, J.A., Weidemann, C.T., McGill, K., Jaggi, J.L., Baltuch, G.H., and Kahana, M.J.  
915 (2009). Human substantia nigra neurons encode unexpected financial rewards. *Science*  
916 323(5920), 1496-1499.

917  
918

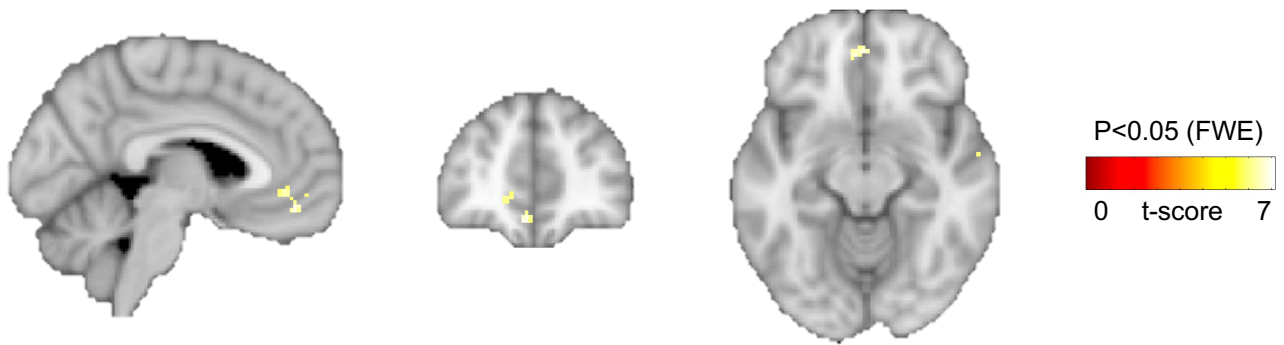
#### 919 **METHODS ONLY REFERENCES**

- 920 Ahn, W.-Y., Haines, N., and Zhang, L. (2017). Revealing neurocomputational mechanisms of  
921 reinforcement learning and decision-making with the hBayesDM package. *Computational*  
922 *Psychiatry* 1, 24-57.
- 923 Andersson, J.L.R., Skare, S., and Ashburner, J. (2003). How to correct susceptibility distortions in spin-  
924 echo echo-planar images: application to diffusion tensor imaging. *Neuroimage* 20(2), 870-888.
- 925 Carpenter, B., Gelman, A., Hoffman, M.D., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M., Guo, J.,  
926 Li, P., and Riddell, A. (2017). Stan: a probabilistic programming language. *Journal of Statistical*  
927 *Software* 76(1), 1-32.
- 928 Gabry, J., and Goodrich, B. (2017). rstanarm: Bayesian applied regression modeling via Stan. R package  
929 version 2.21.1. Retrieved from <https://mc-stan.org/rstanarm/>.
- 930 Gronau, Q.F., Sarafoglou, A., Matzke, D., Ly, A., Boehm, U., Marsman, M., Leslie, D.S., Forster, J.J.,  
931 Wagenmakers, E.-J., and Steingroever, H. (2017). A tutorial on bridge sampling. *Journal of*  
932 *Mathematical Psychology* 81, 80-97.
- 933 McKay, D.J. (2003). *Information Theory, Inference, and Learning Algorithms*. (Cambridge University  
934 Press).
- 935 Papaspiliopoulos, O., Roberts, G.O., and Sköld, M. (2007). A general framework for parameterization of  
936 hierarchical models. *Statistical Science* 22(1), 59-73.
- 937 Vehtari, A., Gelman, A., and Gabry, J. (2017). Practical Bayesian model comparison using leave-one-out  
938 cross-validation and WAIC. *Statistics and Computing* 27(5), 1413-1432.
- 939 Watkins, C.J.C.H., and Dayan, P. (1992). Q-learning. *Machine Learning* 8(3), 279–292.

940  
941  
942



**Supplemental Figure 1. VPRL leads to more accurate learned values of punishing options on the PRP task compared to TDRL.** We used each participant's estimated parameters for the TDRL and VPRL models to compute the expected state-action value (Q-value) for each option on the PRPT task over time. The PRP options are arranged top to bottom as the (A) 25%, (B) 50%, and (C) 75% reward-associated (left column) or punishment-associated (right column) options. Bold green and blue traces represent mean expected value for TDRL and VPRL, respectively, and the shaded region around the means represents one standard error of the mean. TDRL and VPRL model-derived learned values for reward-associated options were very similar to each other, whereas learned values for punishment-associated options were significantly different between models, according to an independent samples t-test at each time point of the difference between the true value (dashed line) and the TDRL or VPRL model-derived learned values across participants. Grey asterisk =  $p < 0.05$ , black asterisk =  $p < 0.05$  Bonferroni corrected.



**Supplemental Figure 2 – Medial prefrontal activity correlates with VPRL-derived expected state-action value signals.** Whole-brain model-based analysis of VPRL learning signals indicates that medial prefrontal cortex parametrically encodes the expected values (VP-Q-value) of the chosen and unchosen options on each trial. Analyses were whole-brain FWE-corrected at  $p < 0.05$ , and the slices in MNI coordinates are  $x = -4$ ,  $y = 46$ ,  $z = -12$ .