

1 Ventral frontostriatal circuitry mediates the computation of reinforcement from
2 symbolic gains and losses

3

4 Hua Tang^{1*}, Ramon Bartolo^{1,2}, Bruno B. Averbeck^{1,3*}

5 1. Laboratory of Neuropsychology, National Institute of Mental Health, National Institutes of Health,
6 Bethesda, MD, USA.

7 2. Laboratory of Sensorimotor Research, National Eye Institute, National Institutes of Health, Bethesda,
8 MD, USA.

9 3. Lead Contact

10 *Correspondence: hua.tang@nih.gov (H.T.), bruno.averbeck@nih.gov (B. B. A.)

11 SUMMARY

12 Reinforcement learning (RL), particularly in primates, is often driven by symbolic outcomes. However, it
13 is usually studied with primary reinforcers. To examine the neural mechanisms underlying learning from
14 symbolic outcomes, we trained monkeys on a task in which they learned to choose options that led to gains
15 of tokens and avoid choosing options that led to losses of tokens. We then recorded simultaneously from
16 the orbitofrontal cortex (OFC), ventral striatum (VS), amygdala (AMY), and the mediodorsal thalamus
17 (MDt). We found that the OFC played a dominant role in coding token outcomes and token prediction
18 errors. The other areas contributed complementary functions with the VS coding appetitive outcomes and
19 the AMY coding the salience of outcomes. The MDt coded actions and relayed information about tokens
20 between the OFC and VS. Thus, OFC leads the process of symbolic reinforcement learning in the ventral
21 frontostriatal circuitry.

22 INTRODUCTION

23 Reinforcement learning (RL) is an adaptive process by which agents learn to make choices to gain rewards
24 over some future time horizon¹. These processes are often studied in animal models using tasks in which
25 choices lead to primary rewards^{2,3}. However, in many situations, choices lead to symbolic outcomes that
26 lead to rewards in the future. Humans are motivated by and will work for symbolic outcomes in the form
27 of money^{4,5}. Animals also readily learn to make choices that lead to symbolic reinforcers, often in the form
28 of tokens. For example, numerous studies have shown that primates will learn to make decisions to
29 maximize the accumulation of tokens that are periodically converted to primary reinforcers⁶⁻⁹.

30 Early studies established an important role for the striatum and the dopamine innervation of the striatum in
31 RL¹⁰⁻¹³. Subsequent work has shown that additional structures, including a network composed of the
32 OFC^{14,15}, AMY^{16,17}, and MDt¹⁸ are also involved in RL¹⁹. Although causal manipulations and
33 neurophysiology have shown that each of these areas plays a role in RL, the results across tasks and areas
34 are not always consistent. For example, lesions to the VS cause deficits in tasks that use probabilistic
35 delivery of primary reinforcers¹⁷. However, these deficits are limited to learning to associate rewards with
36 objects and not actions²⁰. This is not consistent with actor-critic models of RL that suggest that VS stores a
37 general state value representation for policy learning^{21,22}. Similarly, AMY and VS appear to play stronger
38 roles in probabilistic bandit tasks but not in token-based RL tasks, in which tokens are used as
39 reinforcers^{1,7,23}. In token-based RL, VS is only involved in learning to discriminate between relative gains
40 and plays no role in learning to discriminate between gains and losses⁷. The AMY, on the other hand,
41 appears to play a limited role in token-based RL²³. Previous neurophysiology studies have shown that token-
42 based learning may engage cortical networks more than subcortical networks^{8,9}, whereas learning based
43 directly on primary rewards may preferentially engage subcortical networks^{24,25}. Beyond these examples,
44 some learning-related behaviors may rely more on inference processes than incremental value updates that
45 characterize RL and may, therefore, tap into different networks^{26,27}. Thus, the way in which the ventral
46 cortico-striatal-thalamo-cortical network, including the OFC, VS, AMY, and MDt, orchestrates RL,
47 particularly with symbolic reinforcers, is unclear.

48 In the present study, we carried out simultaneous neurophysiology recordings across the ventral cortico-
49 striatal-thalamo-cortical network using a token-based RL task. We examined single neuron, population, and
50 network computations underlying the performance of the task. The results show that the OFC played a
51 dominant role in the computations relevant to learning from symbolic reinforcers. Token outcomes were
52 coded earlier and more strongly in OFC. Conversely, the VS was characterized by specific and unique
53 coding of appetitive choices and outcomes. The AMY was characterized by coding of salience, a finding

54 not apparent in previous work that used only probabilistic appetitive outcomes. The MDt did not appear to
55 contribute a unique computation. However, it played an important role in mediating the interaction of OFC
56 and VS during the calculation of token outcomes. Together, these results define the unique and shared
57 contributions of the ventral frontostriatal circuitry to learning from symbolic reinforcers.

58 RESULTS

59 Two rhesus monkeys were trained on a two-armed bandit task. In this task, the monkeys collected tokens,
60 which were periodically exchanged for juice rewards (Figures 1A-B). In every block of 108 trials, we
61 introduced four novel images. In each trial, two of the four images were presented on the screen. The choice
62 of an image led stochastically to one of -2, -1, +1, or +2 tokens (Figure 1B). The image-outcome
63 relationships were unknown to the monkeys at the start of the block. Monkeys had to learn the values of
64 the images by choosing one of them and observing the outcome. Token outcomes were stochastic such that
65 in 75% of the trials, the monkeys received the number of tokens associated with the chosen image, and in
66 25% of the trials, the number of tokens did not change. Tokens were accumulated across trials and were
67 cashed out every four to six trials, with one drop of juice for each token.

68 The primary reward was apple juice (#juice), but it was only delivered every four to six trials. In each trial,
69 the monkeys had a chance to gain or lose tokens by choosing one of the images. The chosen value
70 information was first carried by the images (cValue), which predicted the change of tokens (Δ token). The
71 number of accumulated tokens (#token) was always presented on the screen. So, value information
72 appeared in different forms, including cValue, Δ token, #token, and #juice (Figure S1A). The following
73 analyses focus on the neural representation and interactions of these forms of value.

74 Choice behavior was influenced by gaining/losing tokens and token numbers

75 We quantified choice behavior by measuring the fraction of times the monkeys chose the image associated
76 with a higher value in each condition (e.g., choose +1 when presented with -1 and +1 images). The monkeys
77 learned to distinguish the values of two images shown on the screen within about 10 trials for the Gain/Loss
78 and Gain/Gain conditions but learned minimally in the Loss/Loss condition (Figure 1G and Figure S1B).
79 We also fit a Rescorla-Wagner (RW) reinforcement learning model to the choice behavior (Figure 1G),
80 which has been explored previously⁷. The model is used below to examine token reward prediction error
81 (RPE). To quantify how the choice behavior (i.e., whether they chose the better option) was affected by
82 task variables, we fit a multi-way ANOVA model to it (Figures 1H-I). The choice behavior was significantly
83 modulated by the number of observations (Trial, $F_{17, 44950} = 109.98$, $p < 0.001$), Gain/Loss condition
84 (Gain/Loss, $F_{1, 44950} = 286.91$, $p < 0.001$), value difference of two options (Δ value, $F_{3, 44950} = 73.37$, $p <$
85 0.001), and the number of tokens before choice (#token, $F_{10, 44950} = 8.68$, $p < 0.001$), but not by the number
86 of trials since last token cashout (CashID, $F_{5, 44950} = 2.07$, $p = 0.066$). The number of observations indicates
87 the learning process, both the Gain/Loss condition and Δ value indicate the updating of tokens, #token

88 indicates the accumulated value, and the CashID indicates the probability of receiving a primary reward.
89 The result indicates that the monkeys understood the meaning of tokens and adjusted their behavior to get
90 more tokens as learning progressed.

91 Neural encoding of choices, primary and symbolic reinforcers

92 We performed simultaneous recordings of population neural activity across four regions of the ventral
93 cortico-striatal-thalamo-cortical network (Figure 1C and Figures S1C-E) from two macaque monkeys using
94 multi-site linear probes. We collected 606 neurons in the 13L region of the orbitofrontal cortex (OFC), 829
95 neurons in the core region of the ventral striatum (VS), 1607 neurons in the basolateral amygdala (AMY),
96 and 1035 neurons in the medial portion of the mediodorsal thalamus (MDt) using a semi-chronic recording
97 procedure (Figures 1D-F). Neuronal activity differed across areas in different task epochs (Figures S1F-G).

98 To examine how choices and rewards were represented in each area, we fit the responses of single neurons
99 with a sliding window ANOVA model. The model included multiple task-relevant factors, which were the
100 number of tokens (#token), the change of token numbers (Δ token), the number of juice drops delivered on
101 cashout trials (#juice), the image pair presented (condition), the stimulus identity (cStim), *a priori* value
102 (cValue, i.e., +2, +1, -1, -2), and direction (cDir) of chosen images. The factors #token, Δ token, and #juice
103 (Figures 2A-C) were the reinforcement signals that drove the monkeys' behavior. Neurons in all areas (>
104 10%), but more in the OFC (> 20%), showed a strong representation of #token (Figures 2A, H). This is
105 consistent with the tokens that were always present on the screen during the trial. There was also a phasic
106 increase in all areas when tokens were updated. The Δ token was also encoded by neurons in all areas but
107 more strongly in the OFC and VS (Figure 2B). OFC also led the encoding of Δ token in time (Figure 2B,
108 inset). The #juice activated more than 40% of neurons in every area (Figures 2C, I). OFC played a
109 significant role in encoding the reinforcement signals, having the highest proportion of neurons encoding
110 these task variables. A substantial proportion of neurons in these areas also showed responses to choices,
111 including the identity (Figure 2D), *a priori* value (Figures 2E, G), and direction (Figure 2F) of chosen
112 images. Learning-related values derived from the RW model were encoded consistently across areas
113 (Figure S4B). The stimulus identities, chosen values, and directions were most strongly represented in the
114 AMY, OFC, and MDt, but differences were often subtle. This result shows that task variables related to the
115 choices and rewards were encoded across the ventral network, with different regions encoding specific
116 types of information at different phases.

117 Diverse encoding of value update information

118 The updating of value information was represented as gaining or losing tokens. In the first analysis, we
119 used only linear encoding (Figure 2B). However, neurons may encode value information in more diverse
120 ways. To quantitatively characterize the encoding patterns across the areas, we examined the activity of
121 each neuron using a series of linear regression models (Figure S2A). We identified the best-fitting model
122 for each neuron and classified them into different functional categories according to how they were tuned
123 to outcomes.

124 Here, we classified neurons encoding $\Delta token$ into five types: neurons encoding (1) value linearly across
125 both Gain and Loss (Figures 3A, F), (2) value salience (Figures 3B, G), (3) categorical Gain/Loss (Figures
126 3C, H), (4) value only for Gain (Figures 3D, I), and (5) value only for Loss (Figures 3E, J). Note that each
127 encoding type could be positively or negatively related to the neural activity. The first category encoded
128 $\Delta token$ linearly (Figure 3A). These neurons encoded the gains and losses on a linear value axis, in other
129 words, processing the gaining and losing of tokens in the same internal value system. Many OFC neurons,
130 but few in the AMY, were of this type. However, more AMY neurons encoded the salience of $\Delta token$.
131 These neurons represent both gains and losses but with inverse correlations of neural activity and value
132 (Figure 3B). The categorical Gain/Loss signal groups each option as gain or loss, regardless of the
133 magnitude. These neurons were found throughout the ventral network, especially in the OFC and VS
134 (Figure 3C). Many more neurons in these areas encoded gains rather than losses. VS and OFC neurons
135 showed robust and phasic responses to gains (Figure 3D). Only a small proportion of neurons encoded
136 losses, mostly in the OFC (Figure 3E). Before the token update, the chosen images predicted the value
137 update. We also classified the neurons encoding $cValue$ into the same five types (Figure S3A). They showed
138 similar encoding patterns, with activity locked to the onset of the images. This analysis shows that these
139 areas played unique roles in encoding the chosen value and token outcome information.

140 An alternative way to examine the encoding of value update information is to measure the representations
141 of gains and losses separately by running two regressions, either using Gain trials (i.e., $\Delta token \in [0, 1, 2]$)
142 or Loss trials (i.e., $\Delta token \in [0, -1, -2]$), then calculate the correlation of the Gain and Loss regression
143 coefficients (Figures 3K-O). The OFC, VS, and MDt populations showed high co-encoding of gains and
144 losses, which indicates the encoding of gains and losses on similar value axes (Figure 3K). Especially for
145 OFC, the effect was consistent for two seconds after the token update (Figure 4L). This was consistent with
146 the single-cell result (Figure 3A), indicating that OFC processes the gaining and losing of tokens closer to
147 objective outcomes than the other areas (Figure 4K; Fisher's z-transformation, $p < 0.05$). Conversely, AMY

148 showed negative correlations between Gain and Loss regression coefficients (Figures 4K, N), indicating
149 population encoding of value salience (Figure 3B).

150 Co-encoding of the value update information

151 In each trial, the monkeys may gain or lose tokens. The value update was predicted by the images (cValue)
152 before the token outcome was delivered (Δ token). Did neurons in these areas encode the two forms of value
153 information in a similar way? To address this question, we first compared the number of neurons that
154 encoded both forms of value linearly. Neurons encoding cValue in a window 500 ms before the token
155 update and Δ token in a window 500 ms after the token update were classified as encoding both cValue and
156 Δ token (blue dots in Figures 4A-D). Fewer neurons in the AMY than in the other areas linearly encoded
157 both factors (Figures 4A-D, F; chi-square test, $p < 0.01$). To address how similar AMY neurons encoded
158 cValue and Δ token, we then calculated the correlation of the cValue and Δ token regression coefficients in
159 each area (Figures S2B-C). AMY neurons showed a significantly lower correlation between cValue and
160 Δ token than the other areas (Figures 4A-D, G; Fisher's z-transformation, $p < 0.05$ for significant neurons,
161 $p < 0.001$ for all neurons). This result shows that neurons in the AMY encoded the value carried by the
162 images and token outcome less consistently, which is also consistent with AMY primarily encoding
163 salience.

164 We also found asymmetric encoding of gains and losses in the VS, OFC and MDt (Figure 4E; chi-square
165 test; VS, $\chi^2 = 95.74$, $p < 0.001$; OFC, $\chi^2 = 19.96$, $p < 0.001$; MDt, $\chi^2 = 10.87$, $p < 0.05$). More VS neurons
166 had positive cValue and Δ token regression coefficients. Together with the results in Figure 3, this indicates
167 that more neurons in the VS fired more when choosing higher-valued images and getting more tokens in
168 the Gain conditions (Figures 4B, E, quadrant 4). OFC neurons showed more balanced encoding in increased
169 and decreased activities (Figure 4E, quadrant 3 vs. 4; chi-square test, $\chi^2 = 3.11$, $p = 0.078$). Thus, neurons
170 in the VS tended to respond with increased firing rates to the choice of good options and outcomes, whereas
171 neurons in the OFC had both positive and negative tuning to the same variables.

172 Co-encoding of primary and symbolic reinforcers

173 Although a growing number of studies have adopted symbolic reinforcers, whether they are encoded in the
174 same way as primary reinforcers remains an open question. To address this, we compared the encoding of
175 #token and #juice. We first calculated the number of neurons encoding each variable. Neurons encoding
176 #token in a window 500 ms after the token update and encoding #juice in a window 500 ms after juice
177 delivery were classified as encoding both #token and #juice (Figure S3B). More neurons in the OFC than

178 in the other areas encoded both #token and #juice (Figure 4H; chi-square test, $p < 0.001$). We then calculated
179 the correlation of the #token and #juice regression coefficients in each area. Surprisingly, neuronal
180 populations in the AMY and VS showed a higher similarity in encoding #token and #juice (Figure 4I and
181 Figure S3C; Fisher's z-transformation, $p < 0.05$). This result shows that the symbolic and primary
182 reinforcers were more similarly encoded in the AMY and VS. Although OFC encoded both at the highest
183 level, it did so with different population responses and, therefore, can discriminate these reinforcers best
184 among the areas.

185 Value updates in the ventral striatum and orbitofrontal cortex

186 We used a stochastic reward schedule such that tokens were only updated in 75% of the trials, and token
187 losses only occurred when tokens could be lost. The encoding of Δ token (Figure 3), therefore, may also
188 include the encoding of RPE (Figure 5A). To address this, we further fit the behaviors with a reinforcement
189 learning model (Figure 1G) and classified each neuron into Δ token or RPE categories, depending on which
190 variable best described the neuron's responses, using linear regression models (Figures S4A-B). RPE was
191 calculated with the RW model. Neurons encoding Δ token responded to the token outcome independently
192 of the learning, whereas neurons encoding RPE encoded the difference between the token outcome and the
193 predicted token outcome, with the prediction estimated by the RW model. Overall, more neurons in the
194 OFC and VS encoded Δ token or RPE than in the other two areas (Figures S4C-F; chi-square test, $p < 0.001$).
195 We then split neurons based on whether their regression coefficients were positive or negative (e.g., whether
196 they increased or decreased their firing rate for RPE; Figures 5B-E). For example, a neuron was called a
197 +RPE neuron when classified in the RPE category and with a positive regression coefficient.

198 Although a similar proportion of neurons in the OFC and VS encoded token outcome information, they
199 showed different patterns. Most neurons in the OFC at the time of token outcome encoded $-\Delta$ token and
200 +RPE (Figure 5B). In other words, they fired more when losing more tokens ($-\Delta$ token) or when they did
201 not lose tokens after choosing a loss option (+RPE). Specifically, this was aligned with the value gradient
202 carried by the choices (Figure 5B inset), even though it was a signal related to the outcome. On the other
203 hand, most neurons in the VS encoded $+\Delta$ token and +RPE (Figure 5C). Neurons in the VS, therefore, fired
204 more for all good outcomes when gaining or not losing more tokens. This was well aligned with the value
205 gradient carried by the outcomes but not the choices (Figure 5C inset). This suggests that value updates
206 were referenced to outcomes in the VS and to choices in the OFC.

207 Population representation of gains and losses

208 Next, we examined how neural populations in each area dynamically encoded gains and losses at the time
209 of choice and token outcome by measuring population response trajectories²⁸. This analysis used pseudo-
210 populations composed of 500 neurons recorded across sessions. We focused on responses in a specific low-
211 dimensional subspace that captured variance due to the eight combinations of cValue and Δ token (Δ token
212 equaled cValue or 0). We first defined the axes of the task-related subspace using linear regression
213 coefficients. Then, the condition-averaged population response was projected on each axis to estimate the
214 representation of the corresponding task variables across time (Figure S2D).

215 The one-dimensional trajectories reflected the dominant tuning at the single-cell level, and correlations
216 between coding of cValue and Δ token at the population level (Figures 3-5). Trajectories diverged into
217 groups aligned with the value gradient following cue onset for cValue (Figures 6A-D) and following token
218 update for Δ token (Figures 6E-H). Because all neurons were z-transformed before these analyses, the axes
219 reflect the strength of encoding the corresponding variable. OFC showed the strongest and most balanced
220 representation of positive and negative choices and outcomes among the four areas, with larger divergence
221 among trajectories in both the cValue (Figure 6A) and Δ token (Figure 6E) axes. The largest deviations for
222 OFC, however, were for Loss conditions. OFC trajectories also showed the population representation of
223 RPEs, with zero token outcomes for Loss options intermediate between Loss and Gain outcomes and zero
224 token outcomes for Gain options merged with Loss outcomes (Figure 6E). Because cValue and Δ token
225 were correlated in OFC, the cValue trajectories also crossed following the token update (Figure 6A).

226 VS population activity also discriminated gains and losses well but with a bias to the Gain conditions. The
227 Gain trajectories had steeper peaks (Figure 6B) and were better separated (Figures 6B, F) than the Loss
228 trajectories. The VS also showed the crossover in cValue trajectories after outcomes (Figure 6B). AMY
229 trajectories were grouped into Gain and Loss groups for cValue (Figure 6C), indicating the overall
230 population coding of Gain vs. Loss. However, the Δ token trajectories, about 250 ms after the token update
231 (Figure 6G), had the weakest responses for all zero token outcomes, and stronger responses for both gains
232 and losses, indicating the overall population coding of the salience of Δ token. MDt trajectories showed
233 similar patterns to those of OFC, including the crossover following zero token outcomes. This suggests that
234 they shared similar information. Overall, these results are consistent with the results shown in Figures 3-5,
235 which confirm the unique contributions of each area in encoding gains and losses.

236 Flow of token information in the ventral network

237 Brain areas function as part of big networks but not as individual isolated areas. The task-relevant
238 information also is not isolated in each area. To understand the flow of token information within the
239 network, we measured the linear dynamics of the trajectories within a trial. This analysis used only
240 simultaneously recorded ensembles and was carried out trial-by-trial. First, we projected the population
241 neural activity from single trials into a 3-D subspace to generate the population state-space representation
242 of #token, Δ token, and cDir (Figure S2D). Then, we put these task-variable-specific latent variables
243 together into a matrix X . The rows included all the latent variables from all the areas, and the columns
244 included each trial (Figure 7A). We stacked all the sessions ($N = 16$) with more than ten neurons recorded
245 from each area simultaneously, then estimated the loading matrix A (Figure 7B). The matrix A
246 characterized the flow of information among variables and areas (Figure S2E), with its columns
247 representing the source areas and variables, and the rows representing the target areas and variables (Figure
248 7B). The matrices were estimated separately at each point in time, as a local linear approximation to the
249 dynamics, which were likely nonlinear.

250 The eigenvalues of the A matrices captured the time constant of the dynamics of information flow among
251 the areas in the network. The top two eigenvalues, which show two peaks, indicate robust information flow
252 during the cue and token update epochs, and the third eigenvalue captured information flow during the
253 token update epoch (Figure 7D). The left and right eigenvectors of the matrix A capture the temporal
254 dynamics of output (Figures S5A-C) and input (Figures S5D-F) information. The first eigenvector also
255 reflected the temporal structure of the task, showing phasic activity that was locked to the associated task
256 variables. The values in the A matrices indicate the strength of information flow between a specific pair of
257 areas and task variables. The temporal dynamics of the information flow varied across the time course of
258 the trial (Figure 7C). We examined the flow of token information across areas and task variables (Figure
259 S6) and found the strongest information flow existed within the same area for most conditions. This is
260 consistent with other work showing that within-region dynamics are higher dimensional than across-region
261 dynamics²⁹.

262 Our dynamics analysis showed good specificity between task variables. For example, despite the fact that
263 cDir, #token, and Δ token were represented across areas, there were minimal interactions between cDir and
264 the value signals (Figures S6F-H). Aiming to address the flow of token information, we focused on the
265 interactions within and between the Δ token and #token. For Δ token (Figure 7E), information flows showed
266 peaks phase-locked to the token update. There was strong reciprocal information flow between the OFC

267 and MDt, reflecting the underlying anatomy¹⁹. We also found flow of information from AMY to the other
268 areas, especially VS. For the #token (Figure 7F), the information flow was relatively continuous along the
269 time course of the trial, which is consistent with the tokens being present on the screen across trials.

270 More interestingly, there was a temporal derivative in the flow of information from #token to Δ token in the
271 OFC and VS (Figure 7G, 100 ms before vs. after derivative; Fisher's z-transformation; VS to VS, $p < 0.001$;
272 OFC to OFC, $p < 0.05$). This shows the mechanism by which Δ token was calculated in the network.
273 Specifically, OFC and VS computed the difference in #token following and prior to the choice (i.e.,
274 $\Delta token = \#token_{after} - \#token_{before}$). Also, we found a stronger flow of information from the #token
275 to Δ token than from the Δ token to #token (Figures S6B, D; Fisher's z-transformation, $p < 0.05$). This
276 suggests that the updating of values in this four-area network originated predominantly in the VS and OFC,
277 consistent with the single-cell results (Figures 5B-C). Because both areas showed time-derivative dynamics,
278 it seems likely that they may both be calculating the update. However, the derivative in the OFC appeared
279 about 100 ms earlier than in the VS, similar to the time of single neurons in the OFC leading Δ token
280 encoding relative to the VS (Figure 2B inset). During this calculation, phasic reciprocal information flow
281 also existed between OFC-MDt and MDt-VS (Figure 7G), indicating that MDt bridged the information
282 flow from the OFC to VS.

283 DISCUSSION

284 We examined the single-cell, population, and network representation of token-based RL across the ventral
285 cortico-striatal-thalamo-cortical network^{1,19,30}. We found that the OFC played an important role in
286 computing information about token updates. While all areas coded this information, OFC coded it earlier
287 and more strongly. We also found that the other areas contributed unique computations. The VS showed a
288 strong bias towards encoding appetitive choices and outcomes, including positive RPEs, with increased
289 firing rates. On the other hand, OFC showed a more balanced encoding of positive and negative choices
290 and outcomes with increased and decreased firing rates. OFC also strongly encoded primary and symbolic
291 reinforcers but distinguished them at the population level. AMY, on the other hand, encoded both
292 reinforcers similarly at the population level. AMY also showed enhanced salience coding and did not show
293 correlations between choice and outcome value coding. Finally, MDt showed enhanced coding of actions
294 and mediated information flow between the OFC and VS during token update calculations.

295 Primary and symbolic reinforcers

296 Human studies typically use symbolic reinforcers, while animal studies use primary reinforcers. However,
297 similarities and differences between the population coding of primary and symbolic reinforcers have not
298 been examined. AMY and, to some extent, VS populations showed stronger correlations between primary
299 and symbolic reinforcers and, therefore, encoded tokens and primary reinforcers similarly. This may
300 explain why earlier lesion work found AMY linked conditioned stimuli to the specific reward properties of
301 the unconditioned stimuli they predicted³¹. OFC showed the most substantial coding of both primary and
302 symbolic reinforcers but encoded them with the lowest similarity, therefore discriminating them well. Thus,
303 OFC represents value information with high fidelity, using a code that preserves detailed information about
304 the type and valence of reinforcement³².

305 Gains and losses

306 We found that VS showed a strong bias towards monotonically encoding gains, specifically coding
307 rewarding choices and outcomes, and better than the expected outcomes, with increased firing rates. Thus,
308 VS consistently signaled token gains with increased firing rates. This is consistent with human imaging
309 studies³³ and the longstanding suggestion that VS is important for motivation^{34,35}. AMY neurons showed
310 enhanced tuning for salience. The AMY population, uniquely among the areas, did not show a correlation
311 between linear choice value and outcome value encoding, and showed negative correlations, at the
312 population level, between gain and loss encoding, both of which are also consistent with salience coding.

313 Previous studies found AMY responses to appetitive and aversive stimuli but could not assess salience
314 because they only used single outcomes for each valence^{3,36,37}. OFC showed monotonic encoding of value
315 across both positive and negative outcomes. Unlike VS and AMY, OFC neurons encoded values with both
316 positive and negative slopes with a slight bias towards neurons responding more for negative value choices
317 and loss outcomes. Recordings in the pregenual anterior cingulate cortex (ACC)³⁸ and the insula⁸ similarly
318 found neuronal populations coding appetitive and aversive choices with both positive and negative slopes.
319 Thus, cortical areas show a more balanced coding of gains and losses than the VS and AMY.

320 OFC contributes to the encoding of symbolic reinforcers

321 We found that OFC encoded accumulated tokens and changes in tokens at both the single-cell and
322 population levels. These findings are consistent with previous work showing that OFC codes state
323 information about the environment^{39,42}, as tokens and token updates define states and state transitions in
324 this task⁴³. State representation is also referred to as a cognitive map, particularly when states have to be
325 inferred^{44,45}. It is a map because states are nodes on graphs, and one has to know the current state and
326 subsequent states to which one can transition. Within reinforcement learning, states are the environmental
327 variables relevant to the learning and action selection process⁴⁶. We also found that OFC encoded choice
328 value, which is consistent with previous work showing that OFC codes economic value⁴⁷. However, we
329 have found that choice value was broadly encoded across our network, and OFC did not encode value more
330 strongly than other areas. We did find differences in the way in which OFC encoded values relative to the
331 other areas, as discussed above. OFC has a high-fidelity representation of the state that includes negative
332 outcomes and the strongest encoding of accumulated tokens.

333 The calculation of value updates

334 We also examined network computations within targeted information dimensions. Analysis of single-trial
335 communication using neurons recorded from multiple regions simultaneously can provide evidence of
336 dynamic network processes that underlie behaviors⁴⁸. Most studies that have examined interactions between
337 areas with neurophysiology data have focused on pairwise interactions between cortical regions and have
338 also not identified dynamic signatures of computations. Instead, they have shown that interactions occur
339 within specific dimensions and often during particular periods within a task. For example, an early study
340 found that choice-relevant signals were relayed from the prefrontal to the parietal cortex to guide behavior⁴⁹.
341 A more recent study found that stronger value coding in the OFC led to accelerated ramping of signals in
342 the ACC⁵⁰. We recently identified object-to-direction information flow among lateral prefrontal cortex

343 subregions in a task in which the location of a valuable stimulus had to be identified before a saccade could
344 be directed toward it to make a choice⁵¹.

345 The current study measured the flow of value information using simultaneous population recordings from
346 four areas in the ventral network. This allowed us to control for multiple, although not all, inputs to each
347 area. We found that token update information was calculated in the OFC and VS as a time-derivative of the
348 information about accumulated tokens (i.e., $\Delta token_{t+1} = \#token_{t+1} - \#token_t$). The signal was earlier
349 in the OFC but stronger in the VS. Thus, the change in accumulated tokens drove information about token
350 outcomes. This is consistent with the finding that OFC and VS strongly encoded $\Delta token$ and token
351 prediction errors at the single-cell level. Note that in previous work with probabilistic reward outcomes, we
352 did not find strong encoding of RPE in these structures^{16,44}. We found minimal support for the alternative
353 hypothesis that token outcomes were integrated to generate information about accumulated tokens (i.e.,
354 $\#token_{t+1} = \#token_t + \Delta token_t$). Thus, analysis of neural dynamics identified an explicit computational
355 correlate of a behavioral process. We also found that MDt mediated interactions between the OFC and VS
356 with respect to token updates by mediating reciprocal information flow between OFC-MDt and MDt-VS.
357 Interestingly, we did not identify a direct interaction between the OFC and the VS, even though they are
358 connected directly^{52,53}. This suggests that the token updates were mediated within the cortical-thalamic-
359 striatal circuitry but not the corticostriatal circuitry.

360 Effects of loss of specific nodes on reinforcement learning

361 These results also provide insight into the task-dependent effects of lesions in previous work. For example,
362 tasks that use probabilistic delivery of primary reinforcers have shown deficits following lesions to all nodes
363 of the ventral network^{14,15,17,18,54,55}. However, lesions of the AMY have almost no effect on learning in the
364 tokens task, and lesions of the VS only affect learning to discriminate gain magnitude but not gains vs.
365 losses^{7,23}. Here, we found that AMY showed enhanced salience coding. Salience coding can be used as a
366 learning signal in the context of probabilistic reward outcomes because salience and gains relative to no
367 outcome are equivalent. However, salience coding cannot be used for learning in token-based reinforcement
368 because large gains and large losses are represented similarly. This may also explain why we found no
369 evidence for the effects of AMY lesions on unsigned prediction error variables in Pearce-Hall models of
370 RL that had been reported previously⁵⁶ when examined in probabilistic reward tasks¹⁷. We also found that
371 VS was strongly biased toward encoding appetitive outcomes. Therefore, lesions to the VS would be
372 expected to have a larger effect on discriminating among gains, as opposed to between gains and losses⁷.
373 Although we did not record in the insula, recent work has shown that the insula has substantial coding of
374 losses in token-based decision-making⁸, consistent with early fMRI results⁵⁷. Therefore, it is possible that

375 manipulations of the insula would show effects specific to loss outcomes. Finally, we found that OFC may
376 preferentially mediate the calculation of token reinforcement, and therefore, manipulations of OFC may
377 have large effects on learning in the tokens task.

378 Conclusion

379 Information about task variables was represented across the ventral network. Although all areas represented
380 task variables, they did so differently. AMY encoded outcome salience and encoded the primary and
381 symbolic reinforcers similarly. VS and OFC encoded value information, with VS strongly biased towards
382 coding positive outcomes with increased firing rates and OFC more balanced towards coding positive and
383 negative outcomes. OFC and VS calculated Δ token as the time-derivative of accumulated tokens, with a
384 shorter latency in the OFC. MDt mediated the interactions about token updates between the OFC and VS
385 during this process. Importantly, the representation and computation of symbolic reinforcement appear to
386 be more strongly mediated by cortical structures, in this case, OFC, than subcortical structures, which may
387 differ from primary reinforcement.

388 ACKNOWLEDGMENTS

389 We thank Christos Constantinidis, Vincent Costa and Diana Burk for their valuable comments, as well as
390 Andy Mitz, Craig Taswell, Miriam Janssen and Sarah Falkovic for their technical help. This work was
391 supported by the Intramural Research Program of the National Institute of Mental Health (ZIA MH002928)
392 to B. B. A. and BBRF NARSAD Young Investigator Award (30892) to H. T. Anatomical MRI scanning
393 was carried out in the Neurophysiology Imaging Facility Core (NIMH, NINDS, NEI).

394 AUTHOR CONTRIBUTIONS

395 Conceptualization, R. B., H. T., and B. B. A.; Methodology, H. T., R. B., and B. B. A.; Investigation, H.
396 T., and R. B.; Visualization, H. T., and B. B. A.; Writing – Original Draft, H. T., and B. B. A.; Writing –
397 Review & Editing, H. T. and B. B. A.; Funding Acquisition, H. T. and B. B. A.; Supervision, B. B. A.

398 DECLARATION OF INTERESTS

399 The authors declare no competing interests.

400 REFERENCES

- 401 1. Averbeck, B., and O'Doherty, J.P. (2022). Reinforcement-learning in fronto-striatal circuits.
402 *Neuropsychopharmacology* 47, 147-162. 10.1038/s41386-021-01108-0.
- 403 2. Jezzini, A., Bromberg-Martin, E.S., Trambaiolli, L.R., Haber, S.N., and Monosov, I.E. (2021). A
404 prefrontal network integrates preferences for advance information about uncertain rewards and
405 punishments. *Neuron* 109, 2339-2352 e2335. 10.1016/j.neuron.2021.05.013.
- 406 3. Pryluk, R., Shohat, Y., Morozov, A., Friedman, D., Taub, A.H., and Paz, R. (2020). Shared yet
407 dissociable neural codes across eye gaze, valence and expectation. *Nature* 586, 95-100.
408 10.1038/s41586-020-2740-8.
- 409 4. Betzel, R.F., and Bassett, D.S. (2017). Multi-scale brain networks. *NeuroImage* 160, 73-83.
410 10.1016/j.neuroimage.2016.11.006.
- 411 5. Chau, B.K., Sallet, J., Papageorgiou, G.K., Noonan, M.P., Bell, A.H., Walton, M.E., and Rushworth,
412 M.F. (2015). Contrasting Roles for Orbitofrontal Cortex and Amygdala in Credit Assignment and
413 Learning in Macaques. *Neuron* 87, 1106-1118. 10.1016/j.neuron.2015.08.018.
- 414 6. Farashahi, S., Azab, H., Hayden, B., and Soltani, A. (2018). On the Flexibility of Basic Risk
415 Attitudes in Monkeys. *J. Neurosci.* 38, 4383-4398. 10.1523/JNEUROSCI.2260-17.2018.
- 416 7. Taswell, C.A., Costa, V.D., Murray, E.A., and Averbeck, B.B. (2018). Ventral striatum's role in
417 learning from gains and losses. *Proc Natl Acad Sci U S A* 115, E12398-E12406.
418 10.1073/pnas.1809833115.
- 419 8. Yang, Y.P., Li, X., and Stuphorn, V. (2022). Primate anterior insular cortex represents economic
420 decision variables proposed by prospect theory. *Nat Commun* 13, 717. 10.1038/s41467-022-
421 28278-9.
- 422 9. Seo, H., and Lee, D. (2009). Behavioral and neural changes after gains and losses of conditioned
423 reinforcers. *J. Neurosci.* 29, 3627-3641. 10.1523/JNEUROSCI.4726-08.2009.
- 424 10. Schultz, W., Dayan, P., and Montague, P.R. (1997). A neural substrate of prediction and reward.
425 *Science* 275, 1593-1599. 10.1126/science.275.5306.1593.
- 426 11. Mohebi, A., Pettibone, J.R., Hamid, A.A., Wong, J.T., Vinson, L.T., Patriarchi, T., Tian, L., Kennedy,
427 R.T., and Berke, J.D. (2019). Dissociable dopamine dynamics for learning and motivation. *Nature*
428 570, 65-70. 10.1038/s41586-019-1235-y.
- 429 12. Steinberg, E.E., Keiflin, R., Boivin, J.R., Witten, I.B., Deisseroth, K., and Janak, P.H. (2013). A
430 causal link between prediction errors, dopamine neurons and learning. *Nat. Neurosci.* 16, 966-
431 973. 10.1038/nn.3413.
- 432 13. Cox, J., and Witten, I.B. (2019). Striatal circuits for reward learning and decision-making. *Nat.*
433 *Rev. Neurosci.* 20, 482-494. 10.1038/s41583-019-0189-2.
- 434 14. Rudebeck, P.H., Saunders, R.C., Lundgren, D.A., and Murray, E.A. (2017). Specialized
435 Representations of Value in the Orbital and Ventrolateral Prefrontal Cortex: Desirability versus
436 Availability of Outcomes. *Neuron* 95, 1208-1220 e1205. 10.1016/j.neuron.2017.07.042.
- 437 15. Rudebeck, P.H., Saunders, R.C., Prescott, A.T., Chau, L.S., and Murray, E.A. (2013). Prefrontal
438 mechanisms of behavioral flexibility, emotion regulation and value updating. *Nat. Neurosci.* 16,
439 1140-1145. 10.1038/nn.3440.
- 440 16. Costa, V.D., Mitz, A.R., and Averbeck, B.B. (2019). Subcortical Substrates of Explore-Exploit
441 Decisions in Primates. *Neuron* 103, 533-545 e535. 10.1016/j.neuron.2019.05.017.
- 442 17. Costa, V.D., Dal Monte, O., Lucas, D.R., Murray, E.A., and Averbeck, B.B. (2016). Amygdala and
443 Ventral Striatum Make Distinct Contributions to Reinforcement Learning. *Neuron* 92, 505-517.
444 10.1016/j.neuron.2016.09.025.

- 445 18. Chakraborty, S., Kolling, N., Walton, M.E., and Mitchell, A.S. (2016). Critical role for the
446 mediodorsal thalamus in permitting rapid reward-guided updating in stochastic reward
447 environments. *Elife* 5. 10.7554/eLife.13588.
- 448 19. Averbeck, B.B., and Murray, E.A. (2020). Hypothalamic Interactions with Large-Scale Neural
449 Circuits Underlying Reinforcement Learning and Motivated Behavior. *Trends Neurosci.* 43, 681-
450 694. 10.1016/j.tins.2020.06.006.
- 451 20. Rothenhoefer, K.M., Costa, V.D., Bartolo, R., Vicario-Feliciano, R., Murray, E.A., and Averbeck,
452 B.B. (2017). Effects of Ventral Striatum Lesions on Stimulus-Based versus Action-Based
453 Reinforcement Learning. *J. Neurosci.* 37, 6902-6914. 10.1523/JNEUROSCI.0631-17.2017.
- 454 21. Joel, D., Niv, Y., and Ruppin, E. (2002). Actor-critic models of the basal ganglia: new anatomical
455 and computational perspectives. *Neural Netw.* 15, 535-547, Pii s0893-6080(02)00047-3.
456 10.1016/s0893-6080(02)00047-3.
- 457 22. Takahashi, Y., Schoenbaum, G., and Niv, Y. (2008). Silencing the critics: understanding the effects
458 of cocaine sensitization on dorsolateral and ventral striatum in the context of an actor/critic
459 model. *Front. Neurosci.* 2, 86-99. 10.3389/neuro.01.014.2008.
- 460 23. Taswell, C.A., Janssen, M., Murray, E.A., and Averbeck, B.B. (2023). The motivational role of the
461 ventral striatum and amygdala in learning from gains and losses. *Behav. Neurosci.* 137, 268-280.
462 10.1037/bne0000558.
- 463 24. Tang, H., Costa, V.D., Bartolo, R., and Averbeck, B.B. (2022). Differential coding of goals and
464 actions in ventral and dorsal corticostriatal circuits during goal-directed behavior. *Cell Rep.* 38,
465 110198. 10.1016/j.celrep.2021.110198.
- 466 25. Parker, N.F., Baidya, A., Cox, J., Haetzel, L.M., Zhukovskaya, A., Murugan, M., Engelhard, B.,
467 Goldman, M.S., and Witten, I.B. (2022). Choice-selective sequences dominate in cortical relative
468 to thalamic inputs to NAc to support reinforcement learning. *Cell Rep.* 39, 110756.
469 10.1016/j.celrep.2022.110756.
- 470 26. Bartolo, R., and Averbeck, B.B. (2020). Prefrontal Cortex Predicts State Switches during Reversal
471 Learning. *Neuron* 106, 1044-1054 e1044. 10.1016/j.neuron.2020.03.024.
- 472 27. Blanco-Pozo, M., Akam, T., and Walton, M.E. (2024). Dopamine-independent effect of rewards
473 on choices through hidden-state inference. *Nat. Neurosci.* 27, 286-297. 10.1038/s41593-023-
474 01542-x.
- 475 28. Mante, V., Sussillo, D., Shenoy, K.V., and Newsome, W.T. (2013). Context-dependent
476 computation by recurrent dynamics in prefrontal cortex. *Nature* 503, 78-84.
477 10.1038/nature12742.
- 478 29. Semedo, J.D., Zandvakili, A., Machens, C.K., Yu, B.M., and Kohn, A. (2019). Cortical Areas Interact
479 through a Communication Subspace. *Neuron* 102, 249-259 e244. 10.1016/j.neuron.2019.01.026.
- 480 30. Averbeck, B.B., and Costa, V.D. (2017). Motivational neural circuits underlying reinforcement
481 learning. *Nat. Neurosci.* 20, 505-512. 10.1038/nn.4506.
- 482 31. Cardinal, R.N., Parkinson, J.A., Hall, J., and Everitt, B.J. (2002). Emotion and motivation: the role
483 of the amygdala, ventral striatum, and prefrontal cortex. *Neurosci Biobehav Rev* 26, 321-352.
484 10.1016/s0149-7634(02)00007-6.
- 485 32. McDannald, M.A., Esber, G.R., Wegener, M.A., Wied, H.M., Liu, T.L., Stalnaker, T.A., Jones, J.L.,
486 Trageser, J., and Schoenbaum, G. (2014). Orbitofrontal neurons acquire responses to 'valueless'
487 Pavlovian cues during unblocking. *Elife* 3, e02653. 10.7554/eLife.02653.
- 488 33. Yacubian, J., Glascher, J., Schroeder, K., Sommer, T., Braus, D.F., and Buchel, C. (2006).
489 Dissociable systems for gain- and loss-related value predictions and errors of prediction in the
490 human brain. *J. Neurosci.* 26, 9530-9537. 10.1523/JNEUROSCI.2915-06.2006.

- 491 34. Mogenson, G.J., Jones, D.L., and Yim, C.Y. (1980). From motivation to action: functional interface
492 between the limbic system and the motor system. *Prog. Neurobiol.* *14*, 69-97. 10.1016/0301-
493 0082(80)90018-0.
- 494 35. Shiflett, M.W., and Balleine, B.W. (2010). At the limbic-motor interface: disconnection of
495 basolateral amygdala from nucleus accumbens core and shell reveals dissociable components of
496 incentive motivation. *Eur. J. Neurosci.* *32*, 1735-1743. 10.1111/j.1460-9568.2010.07439.x.
- 497 36. Belova, M.A., Paton, J.J., Morrison, S.E., and Salzman, C.D. (2007). Expectation modulates neural
498 responses to pleasant and aversive stimuli in primate amygdala. *Neuron* *55*, 970-984.
499 10.1016/j.neuron.2007.08.004.
- 500 37. Paton, J.J., Belova, M.A., Morrison, S.E., and Salzman, C.D. (2006). The primate amygdala
501 represents the positive and negative value of visual stimuli during learning. *Nature* *439*, 865-
502 870. 10.1038/nature04490.
- 503 38. Amemori, K., and Graybiel, A.M. (2012). Localized microstimulation of primate pregenual
504 cingulate cortex induces negative decision-making. *Nat. Neurosci.* *15*, 776-785.
505 10.1038/nn.3088.
- 506 39. Schuck, N.W., Cai, M.B., Wilson, R.C., and Niv, Y. (2016). Human Orbitofrontal Cortex Represents
507 a Cognitive Map of State Space. *Neuron* *91*, 1402-1412. 10.1016/j.neuron.2016.08.019.
- 508 40. Basu, R., Gebauer, R., Herfurth, T., Kolb, S., Golipour, Z., Tchumatchenko, T., and Ito, H.T. (2021).
509 The orbitofrontal cortex maps future navigational goals. *Nature* *599*, 449-452. 10.1038/s41586-
510 021-04042-9.
- 511 41. Costa, K.M., Scholz, R., Lloyd, K., Moreno-Castilla, P., Gardner, M.P.H., Dayan, P., and
512 Schoenbaum, G. (2023). The role of the lateral orbitofrontal cortex in creating cognitive maps.
513 *Nat. Neurosci.* *26*, 107-115. 10.1038/s41593-022-01216-0.
- 514 42. Gardner, M.P.H., and Schoenbaum, G. (2021). The orbitofrontal cartographer. *Behav. Neurosci.*
515 *135*, 267-276. 10.1037/bne0000463.
- 516 43. Burk, D.C., Taswell, C., Tang, H., and Averbeck, B.B. (2024). Computational mechanisms
517 underlying motivation to earn symbolic reinforcers. *J. Neurosci.*, e1873232024.
518 10.1523/JNEUROSCI.1873-23.2024.
- 519 44. Costa, V.D., and Averbeck, B.B. (2020). Primate Orbitofrontal Cortex Codes Information Relevant
520 for Managing Explore-Exploit Tradeoffs. *J. Neurosci.* *40*, 2553-2561. 10.1523/JNEUROSCI.2355-
521 19.2020.
- 522 45. Wilson, R.C., Takahashi, Y.K., Schoenbaum, G., and Niv, Y. (2014). Orbitofrontal cortex as a
523 cognitive map of task space. *Neuron* *81*, 267-279. 10.1016/j.neuron.2013.11.005.
- 524 46. Averbeck, B.B. (2015). Theory of choice in bandit, information sampling and foraging tasks. *PLoS*
525 *Comput. Biol.* *11*, e1004164, e1004164. 10.1371/journal.pcbi.1004164.
- 526 47. Padoa-Schioppa, C., and Cai, X. (2011). The orbitofrontal cortex and the computation of
527 subjective value: consolidated concepts and new perspectives. *Ann. N. Y. Acad. Sci.* *1239*, 130-
528 137. 10.1111/j.1749-6632.2011.06262.x.
- 529 48. Vinck, M., Uran, C., Spyropoulos, G., Onorato, I., Broggin, A.C., Schneider, M., and Canales-
530 Johnson, A. (2023). Principles of large-scale neural interactions. *Neuron* *111*, 987-1002.
531 10.1016/j.neuron.2023.03.015.
- 532 49. Crowe, D.A., Goodwin, S.J., Blackman, R.K., Sakellaridi, S., Sponheim, S.R., MacDonald, A.W., 3rd,
533 and Chafee, M.V. (2013). Prefrontal neurons transmit signals to parietal neurons that reflect
534 executive control of cognition. *Nat. Neurosci.* *16*, 1484-1491. 10.1038/nn.3509.
- 535 50. Balewski, Z.Z., Elston, T.W., Knudsen, E.B., and Wallis, J.D. (2023). Value dynamics affect choice
536 preparation during decision-making. *Nat. Neurosci.* *26*, 1575-1583. 10.1038/s41593-023-01407-
537 3.

- 538 51. Tang, H., Bartolo, R., and Averbeck, B.B. (2021). Reward-related choices determine information
539 timing and flow across macaque lateral prefrontal cortex. *Nat Commun* 12, 894.
540 10.1038/s41467-021-20943-9.
- 541 52. Haber, S.N., and Knutson, B. (2010). The reward circuit: linking primate anatomy and human
542 imaging. *Neuropsychopharmacology* 35, 4-26. 10.1038/npp.2009.129.
- 543 53. Averbeck, B.B., Lehman, J., Jacobson, M., and Haber, S.N. (2014). Estimates of projection overlap
544 and zones of convergence within frontal-striatal circuits. *J. Neurosci.* 34, 9497-9505.
545 10.1523/JNEUROSCI.5806-12.2014.
- 546 54. Walton, M.E., Behrens, T.E., Buckley, M.J., Rudebeck, P.H., and Rushworth, M.F. (2010).
547 Separable learning systems in the macaque brain and the role of orbitofrontal cortex in
548 contingent learning. *Neuron* 65, 927-939. 10.1016/j.neuron.2010.02.027.
- 549 55. Rudebeck, P.H., Ripple, J.A., Mitz, A.R., Averbeck, B.B., and Murray, E.A. (2017). Amygdala
550 Contributions to Stimulus-Reward Encoding in the Macaque Medial and Orbital Frontal Cortex
551 during Learning. *J. Neurosci.* 37, 2186-2202. 10.1523/JNEUROSCI.0933-16.2017.
- 552 56. Li, J., Schiller, D., Schoenbaum, G., Phelps, E.A., and Daw, N.D. (2011). Differential roles of
553 human striatum and amygdala in associative learning. *Nat. Neurosci.* 14, 1250-1252.
554 10.1038/nn.2904.
- 555 57. Pessiglione, M., Seymour, B., Flandin, G., Dolan, R.J., and Frith, C.D. (2006). Dopamine-
556 dependent prediction errors underpin reward-seeking behaviour in humans. *Nature* 442, 1042-
557 1045. 10.1038/nature05051.
- 558 58. Chaure, F.J., Rey, H.G., and Quiñero, R. (2018). A novel and fully automatic spike-sorting
559 implementation with variable number of features. *J. Neurophysiol.* 120, 1859-1871.
560 10.1152/jn.00339.2018.
- 561 59. Hwang, J., Mitz, A.R., and Murray, E.A. (2019). NIMH MonkeyLogic: Behavioral control and data
562 acquisition in MATLAB. *J. Neurosci. Methods* 323, 13-21. 10.1016/j.jneumeth.2019.05.002.
- 563 60. Olejnik, S., and Algina, J. (2000). Measures of Effect Size for Comparative Studies: Applications,
564 Interpretations, and Limitations. *Contemp. Educ. Psychol.* 25, 241-286. 10.1006/ceps.2000.1040.
- 565

566 STAR METHODS

567 KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Deposited data		
Neural data	This paper	DOI
Behavioral data	This paper	DOI
Software and algorithms		
MATLAB R2020b	MathWorks Inc.	https://www.mathworks.com/products/matlab.html
Analysis code	This paper	DOI
MonkeyLogic	NIMH/NIH	https://monkeylogic.nimh.nih.gov/index.html
Adobe Illustrator 2024	Adobe	https://www.adobe.com/products/illustrator.html
Wave_clus 3	Chaure et al., 2018 ⁵⁸	https://github.com/csn-le/wave_clus
Other		
V-probe	Plexon	https://plexon.com/product-category/v-probes/

568 RESOURCE AVAILABILITY

569 Lead contact

570 Further information and requests for resources and reagents should be directed to and will be fulfilled by
571 the lead contact, Bruno B. Averbeck (bruno.averbeck@nih.gov).

572 Materials availability

573 This study did not generate new unique reagents.

574 Data and code availability

575 The datasets supporting the current study will be publicly available as of the date of publication. DOIs are
576 listed in the key resources table.

577 All original codes will be publicly available as of the date of publication. DOIs are listed in the key resources
578 table.

579 Any additional information required to reanalyze the data reported in this paper is available from the lead
580 contact upon request.

581 EXPERIMENTAL MODEL AND STUDY PARTICIPANT DETAILS

582 Subjects

583 The experiments were performed on one adult male (10 kg) and one adult female (7.5 kg) rhesus macaque
584 (*Macca mulatta*). They were 8-10 years old. The monkeys were pair-housed when possible and had access
585 to food 24 hours per day. On testing days, the monkeys were placed on water control and earned their juice
586 through performing the task. On non-testing days, the monkeys were given *ad libitum* access to water.

587 Ethics

588 Experimental procedures for all monkeys followed *the Guide for the Care and Use of Laboratory Animals*
589 and were approved by the National Institute of Mental Health Animal Care and Use Committee.

590 METHOD DETAILS

591 Experimental Setup

592 Monkeys were trained to perform a saccade-based two-armed bandit task. Stimuli were presented on a 19-
593 inch LCD monitor situated 40 cm from the monkeys' eyes. During training and testing, the monkeys sat in
594 a primate chair with their heads restrained. Stimulus presentation and behavioral monitoring were
595 controlled by MonkeyLogic⁵⁹. The eye movements were monitored at 400 fps using a Viewpoint eye tracker
596 (Arrington Research, Scottsdale, AZ) and sampled at 1 kHz. A fixed amount of apple juice was delivered
597 through a pressurized plastic tube gated by a solenoid valve on rewarded trials.

598 Task Design

599 The task was developed and first used in our previous study ⁷. Each session had nine blocks. Each block
600 used four novel images associated with different values, including +2, +1, -1, and -2. The monkeys
601 obtained more tokens by choosing images associated with larger values. Choosing one of the images led to
602 gaining or losing a corresponding number of tokens. Also, token numbers could not be negative, so
603 choosing a loss image when there were no accumulated tokens had no effect. We used a stochastic reward
604 schedule. The number of tokens updated 75% of the time and did not change 25% of the time. To complete
605 a trial successfully, the monkey first acquired and held central fixation for 400-600 ms. Then, two of the
606 four images were randomly selected by the computer and displayed on the screen. The animal made its
607 selection by saccading to one of them. The unchosen image disappeared when the monkey reached the
608 target. Saccade fixation was maintained on the chosen image for 500 ms. After that, the chosen image also
609 disappeared, and the token number was updated according to the chosen image. Tokens accumulated across

610 trials and were cashed out for juice every four to six trials, with the interval randomly selected. At cash-
611 out, the animals were given one drop of juice for each token. When each drop of juice was delivered, one
612 token was removed from the screen. Up to 12 tokens were accumulated and were displayed on the screen
613 across trials.

614 The task had six individual conditions defined by the possible values of the image pairs. The conditions
615 within a block of 108 trials (6 conditions \times 2 counterbalanced for left and right sides \times 9 repetitions) were
616 presented pseudo-randomly. The animals saw each condition twice, once on the left and once on the right,
617 every 12 trials before seeing any condition a third time. We introduced four novel images at the beginning
618 of each block. Images provided as choice options were normalized for luminance and spatial frequency
619 using the SHINE toolbox for MATLAB, described previously⁷.

620 Surgical procedures

621 Each monkey was surgically implanted with a titanium headpost, and a 25 \times 35 mm recording chamber to
622 allow vertical grid access to the OFC, VS, AMY, and MDt (Figures 1D-F). Grid holes for the MDt had 16°
623 angles to allow better access to the target area (Figure S1D). Chamber placements were planned and verified
624 with T1 and T2 magnetic resonance imaging (MRI, 3.0 T). Small burr holes were drilled above each target
625 area. A grid was installed, and one guide tube was inserted through each burr hole. The guide tubes were
626 lowered to about 1-4 mm (depending on the target areas) above the target areas and were glued to the grid.
627 We removed the guide tubes and placed new ones in the adjacent locations after 3-5 recording days (Figures
628 S1D-E). The locations of the guide tubes were verified with MRI after guide tube replacement. All sterile
629 surgeries and MRI scans were performed under anesthesia.

630 Neurophysiological Recordings

631 Neurophysiology recordings (Figures S1C-D) began after the monkeys had recovered from the surgery. We
632 lowered one linear electrode array (V-probe, Plexon Inc, Dallas, TX) into each guide tube on every
633 recording day. Thirty-two channel electrodes with 150 μ m inter-contact spacing probes were used in the
634 VS and MDt, and 64-channel electrodes with 150 μ m inter-contact spacing probes were used in the OFC
635 and AMY. The probes were advanced to their target location by a four-channel micromanipulator (NAN
636 Instruments, Nazareth, Israel) attached to the recording chamber. The depths of the neurons were estimated
637 by their recording locations relative to the tip of the guide tubes (verified with MRI). Electrophysiological
638 data were acquired with a 512-channel Grapevine System (Ripple, Salt Lake City, UT). The spike
639 acquisition threshold was set at a 4.0 \times root mean square (RMS) of the baseline signal for each electrode.
640 Behavioral event markers from MonkeyLogic and eye-tracking signals from Viewpoint were sent to the

641 Ripple acquisition system. The extracellular signals were high-pass filtered (1 kHz cutoff) and digitized at
642 30 kHz to acquire the single-cell activity. Spikes were sorted offline via Wave_clus 3⁵⁸.

643 Choice behavior

644 Each block had six conditions. The conditions pseudo-randomly appeared in the task. The number of
645 observations in each condition increased as a function of learning. We quantified choice behavior during
646 the task by measuring the fraction of choosing the image associated with a higher value in each condition.
647 We then measured how different task variables affected the choice behavior by fitting a multi-way ANOVA
648 model. Factors, including the number of observations (Trial), Gain/Loss condition (Loss/Loss, Loss/Gain,
649 or Gain/Gain), value difference of the options (Δ value, e.g., Δ value = 4 in the condition of -2 vs. +2), the
650 number of tokens before choice (#token), and the number of trials since last token cashout (CashID) were
651 used in the model. The model was run session by session to measure the contribution of each variable in
652 each session, using all the sessions to acquire the statistics for each variable. All trials in which monkeys
653 chose one of the two stimuli were analyzed. Trials in which the monkey broke fixation, failed to make a
654 choice, or attempted to saccade to more than one target were excluded.

655 Effect size

656 To quantify the contribution of each task variable to the monkey's choice. We computed each factor's effect
657 size, ω^2 , from the ANOVA model output. ω^2 is an unbiased estimator of the amount of variance in neural
658 activity explained by each task variable, and ranges between -1 and 1⁶⁰. It is given by:

$$659 \quad \omega^2 = \frac{df_{effect} \times (MS_{effect} - MS_{error})}{SS_{total} + MS_{error}}$$

660 where df_{effect} refers to the degrees of freedom associated with the factor, MS_{effect} refers to the mean
661 squares, MS_{error} is the mean squared error, SS_{total} is the sum of squares of all factors.

662 Responsive neurons

663 To identify neural responses to different task components, we fit a sliding window multi-way ANOVA
664 model to spike counts computed in 200 ms bins, advanced in 50 ms increments, and time-locked to token
665 update. Factors, including the number of tokens on the screen (#token, may change after the choices), the
666 change of token numbers (Δ token), the drops of juice delivered (#juice), image pairs with different value
667 combinations presented (condition, e.g., +2 vs. -1), the order of blocks (blockID), the identity (cStim), a
668 *priori* value (cValue), and direction (cDir) of chosen images were used in the model. The stimulus identity

669 was the specific image used in each block to represent each outcome, which was the interaction of cValue
670 and blockID. The factor blockID was used to remove non-stationarity due to drift.

671 **Statistic test for the proportion of neurons**

672 The binomial test was applied to test whether the proportion of responsive neurons was significantly above
673 the chance level (5% most of the time). The chi-square test was used to compare the proportions of
674 responsive neurons between different pairs of brain areas or among four areas. Significant encoding for
675 each factor at each time bin was evaluated at $p < 0.05$. A neuron that showed a significant response to a
676 factor in no less than three contiguous bins in the statistics was considered to be responsive to that factor.

677 **Neuronal coding regression analysis**

678 To quantitatively characterize how the updating of values (including cValue and $\Delta token$) was encoded in
679 different areas, we examined the activity of each neuron using a series of multivariate linear regression
680 models⁸. The dependent variable, neural activity, was first z-scored by subtracting the mean response from
681 the firing rate at each time and in each trial and dividing the result by the standard deviation of the responses.
682 Both the mean and the standard deviation were computed by combining the neurons' responses across all
683 trials and times. We then described the z-scored responses of neuron i at time t as a linear combination of
684 several task variables. The independent variables were the same as the factors used in the ANOVA model:

$$685 \quad r_{i,t}(k) = \beta_{i,t,cValue} * cValue(k) + \beta_{i,t,\Delta token} * \Delta token(k) + \dots + \beta_{i,t,cDir} * cDir(k) + \varepsilon$$

686 where $r_{i,t}(k)$ is the z-scored response of neuron i at time t and on trial k , $\Delta token(k)$ is the change of
687 tokens on trial k . The regression coefficients, $\beta_{i,t,f}$, describe how much the trial-by-trial firing rate of
688 neuron i , at a given time t during the trial, depends on the corresponding task variable f .

689 We tested all potential combinations of tuning for $\Delta token$ (six forms: General value, Saliency, Gain/Loss,
690 Gain, Loss, nan) and cValue (six forms: General value, Saliency, Gain/Loss, Gain, Loss, nan). Nan means
691 the variable was removed from the model. A particular variable could only be represented by one specific
692 form but not by combinations of more than one form (e.g., a model containing the $\Delta token$ variable could
693 include Gain or Saliency but not both). In total, we tested $6 \times 6 = 36$ models for each neuron. This method
694 is also illustrated in Figure S2A. We determined the best-fitting model for each neuron using the Akaike
695 information criterion (AIC). Neurons were classified into different functional categories (up to one category
696 per variable) according to the combination of the forms of these two variables included in the best-fitting
697 model. We computed the AIC:

698
$$AIC = 2p - 2 \ln(\hat{L})$$

699 where p is the number of free parameters in the model, and \hat{L} is the maximized value of the likelihood
700 function.

701 Population coding similarity

702 The regression coefficient, $\beta_{i,t,f}$, from the multivariate linear regression model reflects the weight of one
703 task variable, f , in explaining the variation of the neuron's activity at time t . The regression coefficients,
704 $\vec{\beta}_f$, from a neural population, represent the weights of each neuron in the population encoding one task
705 variable. Hence, computing the correlation between two regression coefficient vectors from the same neural
706 population tells us how similar a neural population encodes two task variables (Figures S2B-C):

707
$$r = \text{corr}(\vec{\beta}_{f1}, \vec{\beta}_{f2})$$

708 where $\vec{\beta}_f$ is a vector that consists of regression coefficients of task variable f in a neural population. To
709 assess the significance of the correlation between two regression coefficients, we transformed the
710 correlation coefficients into normally distributed z-scores using Fisher's z-transformation.

711 Reward prediction error

712 The RPE was defined as the difference between the change of tokens $\Delta token(t)$ and the estimated value
713 of the chosen image, which is given by:

714
$$RPE(t) = \Delta token(t) - v_i(t)$$

715 where t means trial order, and i means the chosen option among two options. The updating of value v_i was
716 estimated using the Rescorla–Wagner equation:

717
$$v_i(t + 1) = v_i(t) + \rho_h RPE(t)$$

718 where ρ is the learning rate. The updated value estimate $v_i(t + 1)$ equals the previous value estimate $v_i(t)$
719 plus the RPE scaled by the cue-dependent learning rate ρ_h for images associated with different *a priori*
720 values h . These values were then passed through a soft-max function to give choice probabilities for the
721 image pairs presented in each trial:

722
$$d_j(t) = (1 + e^{\beta(v_i(t) - v_j(t))})^{-1}, d_i(t) = 1 - d_j(t)$$

723 where β is the choice consistency or inverse temperature parameter, fit across all six cue conditions, and i
724 and j are the two choice options. We then maximized the likelihood of the animal's choices, D , given the
725 parameters, using the cost function:

$$726 \quad f(D|\alpha_i, \beta) = \prod_t [d_1(t)c_1(t) + d_2(t)c_2(t)]$$

727 where $d_1(t)$ is the choice probability value for one option on trial t , and $c_1(t)$ and $c_2(t)$ are indicator
728 variables that take on a value of 1 if the corresponding option was chosen, and 0 otherwise. This model was
729 fit across blocks in each session for each monkey to give one set of free parameters for each session.

730 Targeted dimensionality reduction

731 We used the regression coefficients described above to identify dimensions in state space representing each
732 task variable²⁸. For each variable, we first build a set of coefficient vectors $\overrightarrow{\beta_{f,t}}$ whose entries $\beta_{f,t}(i)$
733 correspond to the regression coefficient for task variable f , time t , and neuron i . The vectors $\overrightarrow{\beta_{f,t}}$ (of length
734 N_{unit}) are obtained by simply rearranging the entries of the vectors $\overrightarrow{\beta_{i,t}}$ (of length N_f) computed above.
735 Each vector, $\overrightarrow{\beta_{f,t}}$, thus corresponds to a direction in state space that accounts for variance in the population
736 response at time t , due to variation in task variable f .

737 We used time-dependent regression vectors, B_t , which is a matrix with each column corresponding to the
738 one regression vector $\overrightarrow{\beta_{f,t}}$. We referred to them as the 'task-related axes'. These axes span a 'regression
739 subspace' representing the task-related information the neural population coded. We then projected the
740 single-trial population responses onto these axes to study the representation of the task-related variables in
741 each trial:

$$742 \quad X_t = B_t^T R_t$$

743 where R_t is single-trial population response at time t , which is of dimension $N_{unit} \times N_{trial}$. And X_t is the
744 set of time-series vectors over all task variables and trials, which is of dimension $N_f \times N_{trial}$. It represents
745 the population coding of task-related information at time t of every single trial. This method is also
746 illustrated in Figure S2D.

747 Linear dynamics

748 We fit an autoregressive, linear model to the X_t ,

749
$$X(t + 1) = AX(t) + \varepsilon(t)$$

750 where $X(t) = [X_{OFC}(t), X_{VS}(t), X_{AMY}(t), X_{MDt}(t)]^T$ is a matrix of dimension $(N_{area} \cdot N_f) \times N_{trial}$, and
751 $X_{area}(t) = [x_{k,\#token}(t), x_{k,\Delta token}(t), x_{k,cDir}(t)]$ is a matrix of dimension $N_{trial} \times N_f$. $\varepsilon(t)$ is the noise
752 term. A is a 12×12 dynamics matrix, where the $12 = 4 \times 3$ dimensions correspond to the combinations of
753 brain areas and task variables. The model was fit using data points in a non-overlapping 25 ms moving
754 window. This resulted in a time-dependent estimate of the matrix A . Eigenvalues of A closer to 0 indicate
755 a faster decay, and eigenvalues near 1 would correspond to a slower decay. The method is also illustrated
756 in Figures 7A-C and Figure S2E.

757 A z-test was applied to test for significant differences between the actual and shuffled results for the
758 elements of the A matrices. The shuffled data were generated using the same matrix X but with randomized
759 trial indexes across rows ($areas \times variables$). The actual results with a mean outside the 99% confidence
760 interval of the shuffled results showed a significant response.

761 One-dimensional task-variable-specific trajectory

762 We also carried out the targeted dimensionality reduction using condition-averaged pseudo-populations
763 composed of 500 neurons recorded across sessions to study the population representation of gains and
764 losses. We projected the condition-averaged population responses onto these axes:

765
$$X_t = B_t^T R_t$$

766 where R_t is condition-averaged population response at time t , which is of dimension $N_{unit} \times N_{condition}$.
767 And X_t is the set of time-series vectors over all task variables and conditions, which is of dimension
768 $N_f \times N_{condition}$.

769 The variability of trajectories across different conditions was estimated using a bootstrap procedure.
770 According to the null hypothesis, we generated data with no differences among the conditions by sampling
771 with replacement from the combined set of all conditions. We then calculated the sum of the standard error
772 of each condition from the mean of all conditions for each bootstrap condition set. We did this 1000 times.
773 It gave us 1000 different sampled standard errors from the null distribution. We then compared the standard
774 error of the actual data to the standard errors in the null distribution. If the actual standard error is in the
775 99% confidence interval of the null distribution, the trajectories significantly differ among the conditions
776 (i.e., $p < 0.01$).

777 QUANTIFICATION AND STATISTICAL ANALYSIS

778 Unless otherwise indicated, all data were presented as means \pm SEM (standard error of the mean). The
779 statistical analyses performed were indicated in the main text and detailed in STAR Methods. Statistical
780 comparisons were analyzed in MATLAB (Mathworks), including parametric, non-parametric, and
781 permutation-based statistics, as detailed in STAR Methods. Figures were prepared with MATLAB.

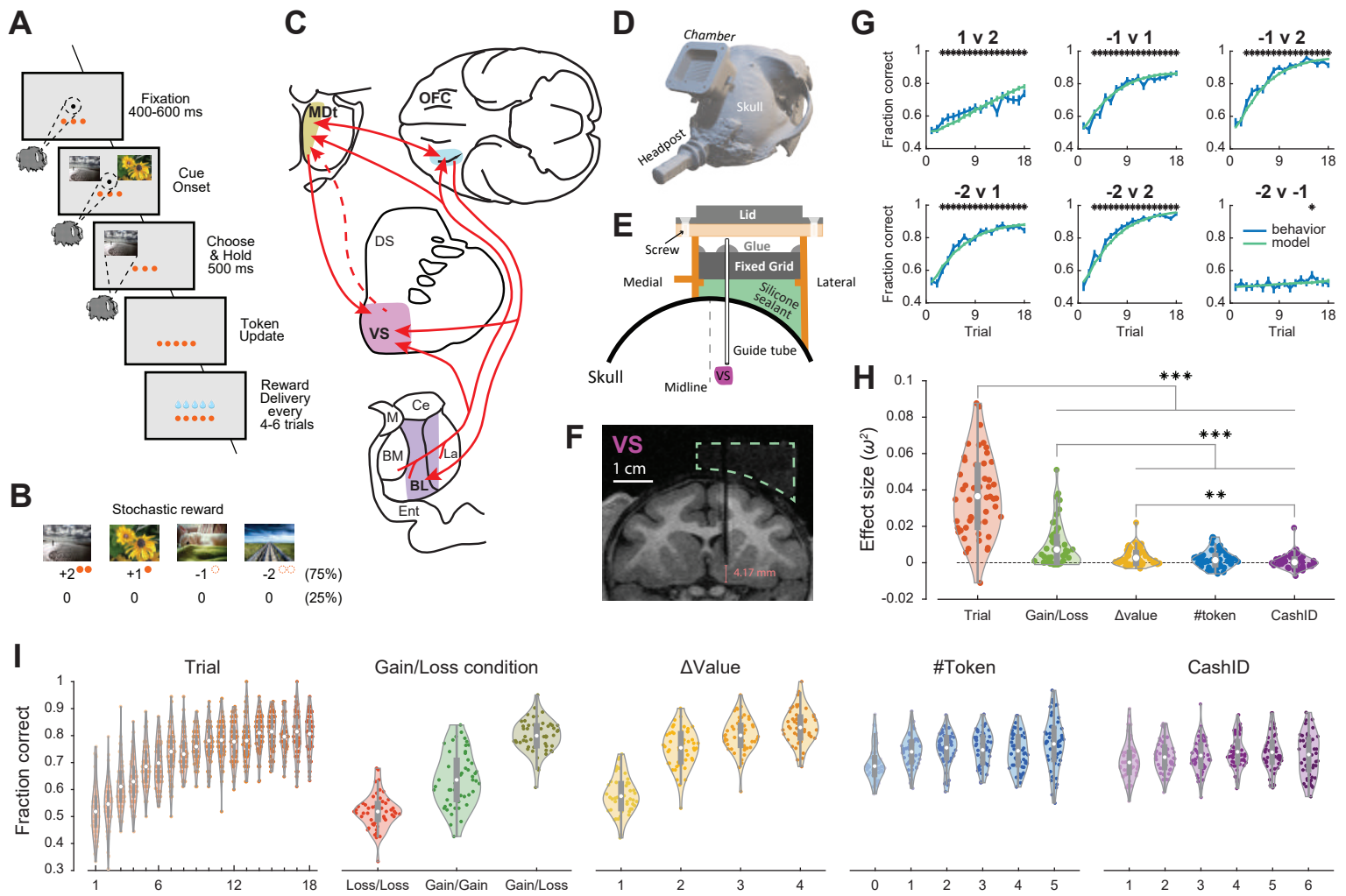


Figure 1

782 FIGURE LEGENDS

783 Figure 1. Task, behavior, recording method, and recorded areas.

784 (A-B) Behavioral task. (A) Structure of an individual trial. Successive frames illustrate the sequence of
785 events. In each trial, only two of the four images were presented. Monkeys chose between them and gained
786 or lost (stochastic: 75% change, 25% no change) the corresponding number of tokens. Choices could be
787 made as soon as the images were shown. Accumulated tokens were cashed out every four to six trials, with
788 one drop of juice for each token. (B) In each block, the monkeys learned the values (+2, +1, -1, -2) of four
789 novel images.

790 (C) Schematic recording areas, including the orbitofrontal cortex (OFC), ventral striatum (VS), basolateral
791 amygdala (AMY), and mediodorsal thalamus (MDt). Red arrows indicate the anatomical connections
792 between these areas. The dashed line indicates an indirect connection.

793 (D-F) Semi-chronic recording. (D) A chamber was implanted. (E) Then, guide tubes were inserted into each
794 target area through a grid. The guide tubes stopped above the target areas and stayed in the brain across
795 recording days. The empty space in the chamber was filled with silicone sealant to prevent potential
796 infections. Multi-site linear probes were lowered into target areas through the guide tubes every recording
797 day. (F) A coronal section view of the chamber shows the shadow of a guide tube and the target area under
798 MRI. The dashed green polygon indicates the silicone sealant.

799 (G-I) Choice behavior. (G) Fraction of choosing the image with a higher value in each condition. Blue lines
800 represent monkey behavior, and green lines represent the prediction of the reinforcement learning model.
801 Error bars represent mean \pm SEM. The black asterisks indicate a significant difference in monkey behavior
802 from the change level (one-sample t-test, $p < 0.01$). (H) The effect size of each task variable. Paired t-test,
803 $**p < 0.01$, $***p < 0.001$. Non-significant pairs are not indicated. (I) Choice behavior as a function of the
804 task variables. Results were averaged from two monkeys, $n = 50$ sessions. Violin plots showing the range
805 of values calculated across all sessions (center dot, mean; box, 25th to 75th percentiles; whiskers, $\pm 1.5 \times$
806 the interquartile range; dots, each session; shade, density curve).

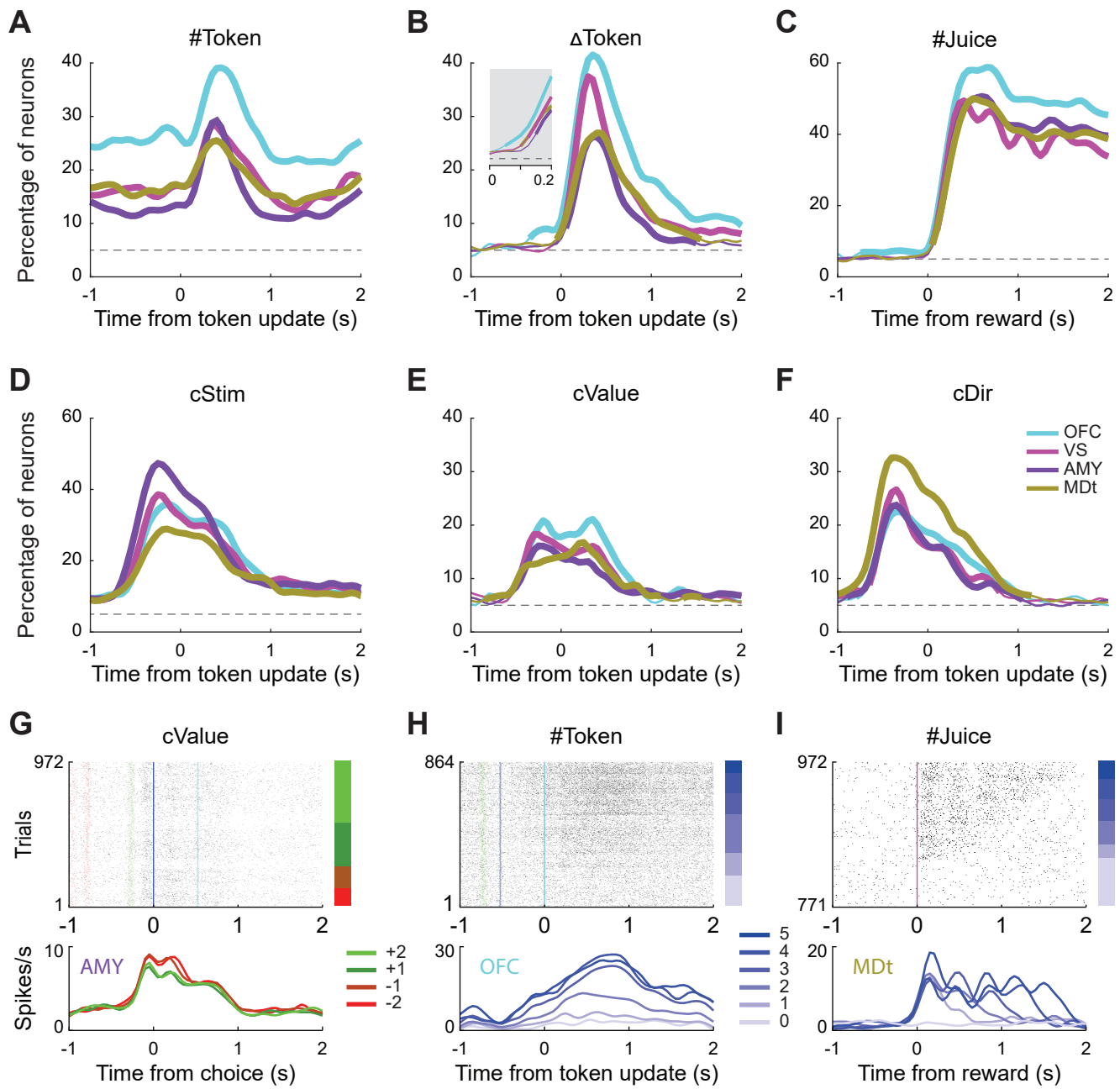


Figure 2

807 **Figure 2. Neural encoding of choices, primary and symbolic reinforcers.**

808 (A-F) Percentage of neurons in each area encoding #token (A), Δ token (B), #juice (C), the identity (D), *a*
809 *priori* value (E), and direction (F) of chosen images. Data were aligned to either the token update or the
810 onset of juice delivery. Inset in B: ANOVA with 25 ms bin showing response latencies to Δ token different
811 among areas. Dashed horizontal lines represent the chance level. Thick lines indicate a significant
812 difference between the corresponding area and chance level (binomial test, $p < 0.01$).

813 (G-I), Example neurons encoded the *a priori* value of the chosen images (G), #token (H), and #juice (I).
814 Each row in the raster plot represents the spikes during a trial. Red, green, blue, cyan, and magenta lines
815 represent fixation, cue onset, choice, token update, and juice delivery. The bars on the right side represent
816 the trial groups under different conditions. The bottom curves indicate mean activities, split by trial groups
817 defined by the bars.

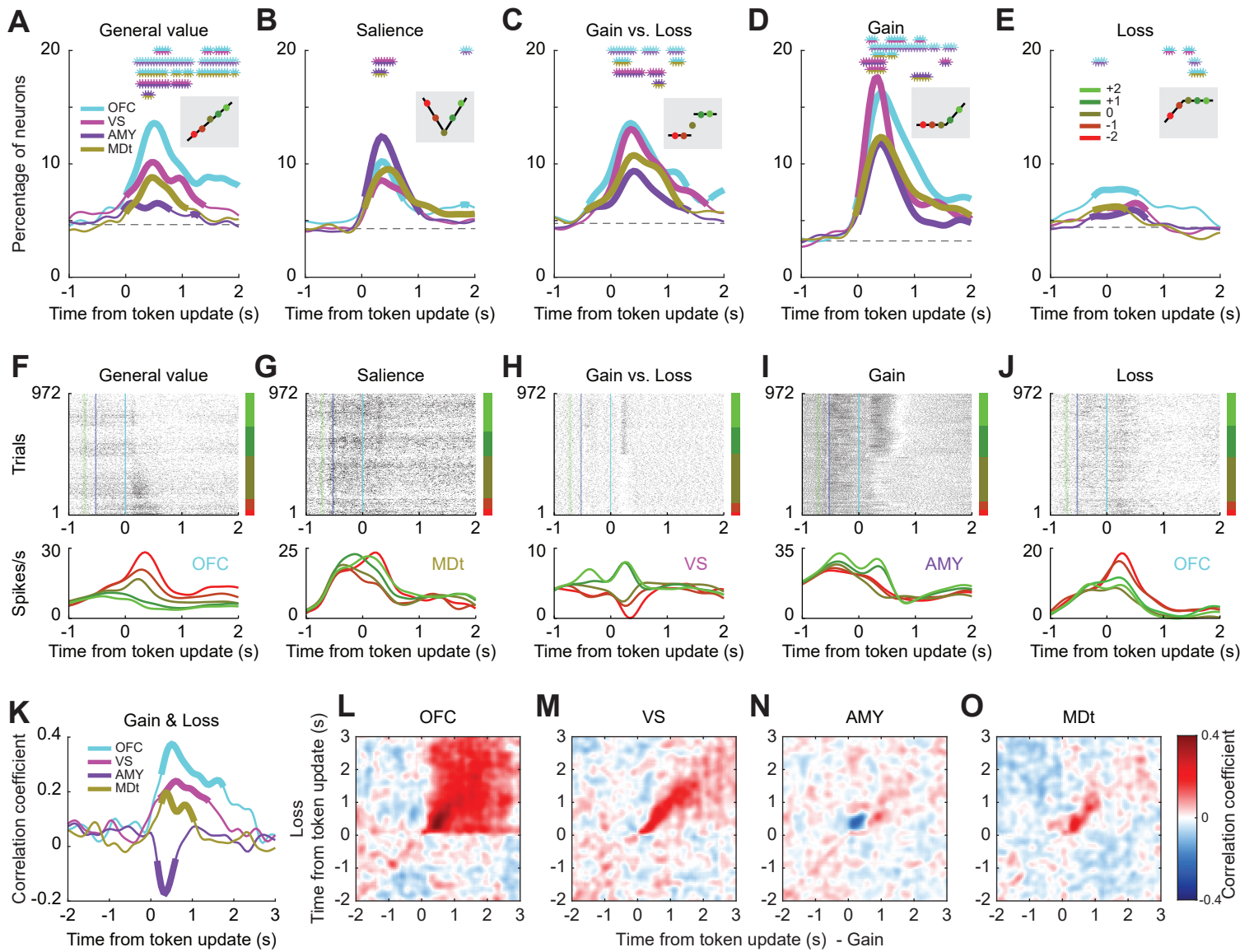


Figure 3

818 **Figure 3. Diverse representation of gains and losses.**

819 (A-E) Percentage of neurons in each area encoding Δ token. The insets indicate the corresponding tuning
820 curves. The coefficients could be positive or negative, meaning the tuning curves could be flipped vertically.

821 (A) General value: monotonic encoding values across gain and loss contexts. (B) Saliency: monotonic
822 encoding values across gain and loss contexts but with opposite directions. (C) Gain/Loss: categorical
823 encoding values in gain and loss contexts. (D) Gain: encoding values only in the gain context. (E) Loss:
824 encoding values only in the loss context. Dashed horizontal lines represent chance levels, defined by the
825 mean of all areas during the fixation epoch. Thick lines indicate a significant difference between the
826 corresponding area and chance level (binomial test, $p < 0.05$). The double-colored asterisks indicate a
827 significant difference between pairs of areas indicated by the colors (chi-square test, $p < 0.05$).

828 (F-J), Example neurons for each category.

829 (K-O) Co-encoding of gains and losses. Correlations of the Gain and Loss regression coefficients at the
830 same (K) or cross (L-O) time points for each area. Thick lines in (K) indicate a significant difference
831 between the corresponding area and baseline level during the fixation period (binomial test, $p < 0.05$).

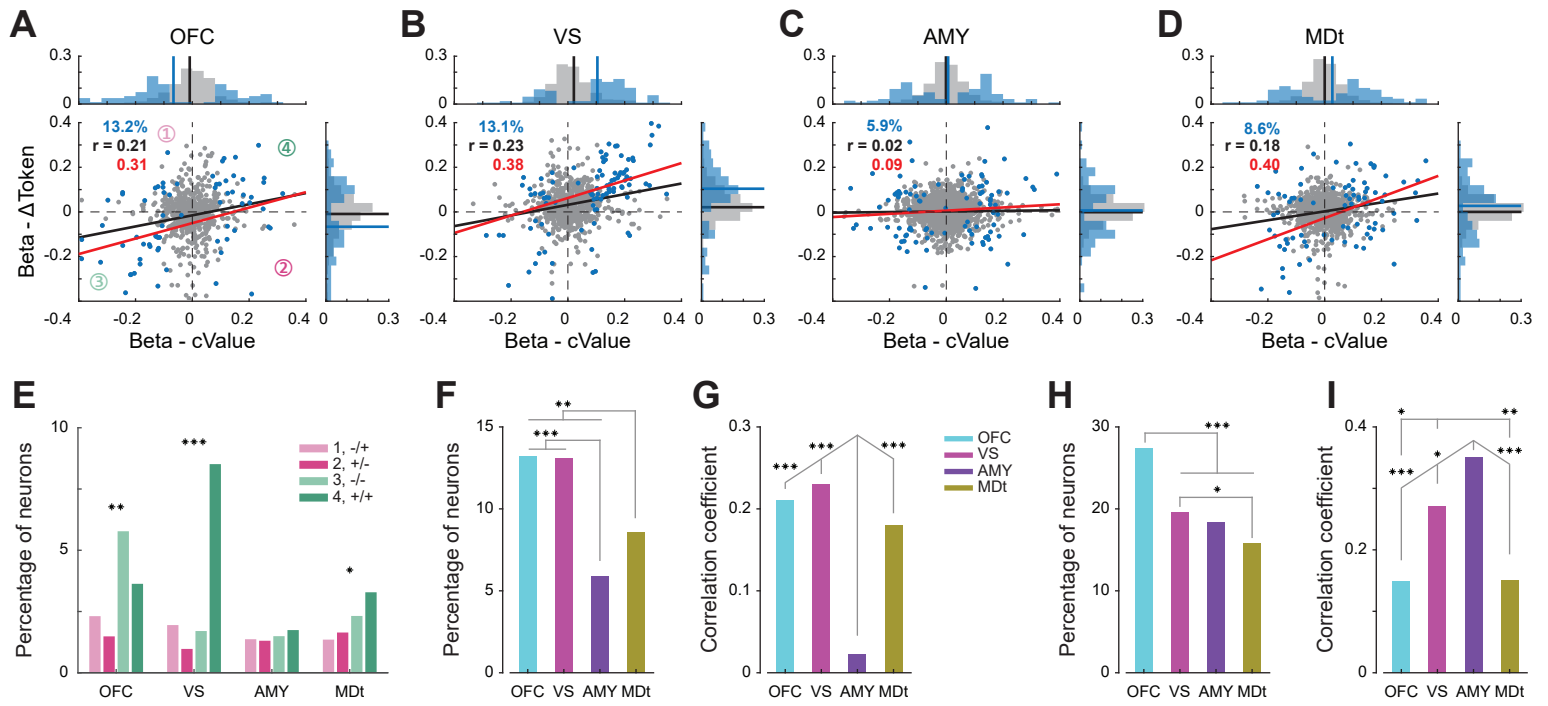


Figure 4

832 **Figure 4. Population encoding of value update and reinforcers.**

833 (A-G) Co-encoding of cValue and Δ token. (A-D) The x- and y-axes represent the regression coefficients
834 of cValue and Δ token. The blue dots represent the neurons encoding both of them. The gray dots represent
835 the other neurons. Percentages on the left corner indicate the proportions of neurons encoding the cValue
836 and Δ token. Black and red lines are the lines of best fit for all the neurons and significant neurons, and r
837 indicates the corresponding correlation coefficient. The bars summarize the distributions of blue and gray
838 dots, with blue and black lines representing the means. (E) Summary of the distributions of blue dots in
839 each quadrant (also indicated in A) for each area. Chi-square test, * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. (F)
840 The proportion of neurons encoding cValue and Δ token in each area. Summary of the percentages in (A-
841 D). Chi-square test, ** $p < 0.01$, *** $p < 0.001$. (G) Encoding similarity between cValue and Δ token in each
842 area. Summary of the black r in (A-D). Fisher's z-transformation, *** $p < 0.001$. Non-significant pairs are
843 not indicated.

844 (H-I) Co-encoding of #token and #juice. (H) The proportion of neurons encoding #token and #juice in each
845 area. Chi-square test, * $p < 0.05$, *** $p < 0.001$. (I) Encoding similarity between #token and #juice in each
846 area. Fisher's z-transformation, * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

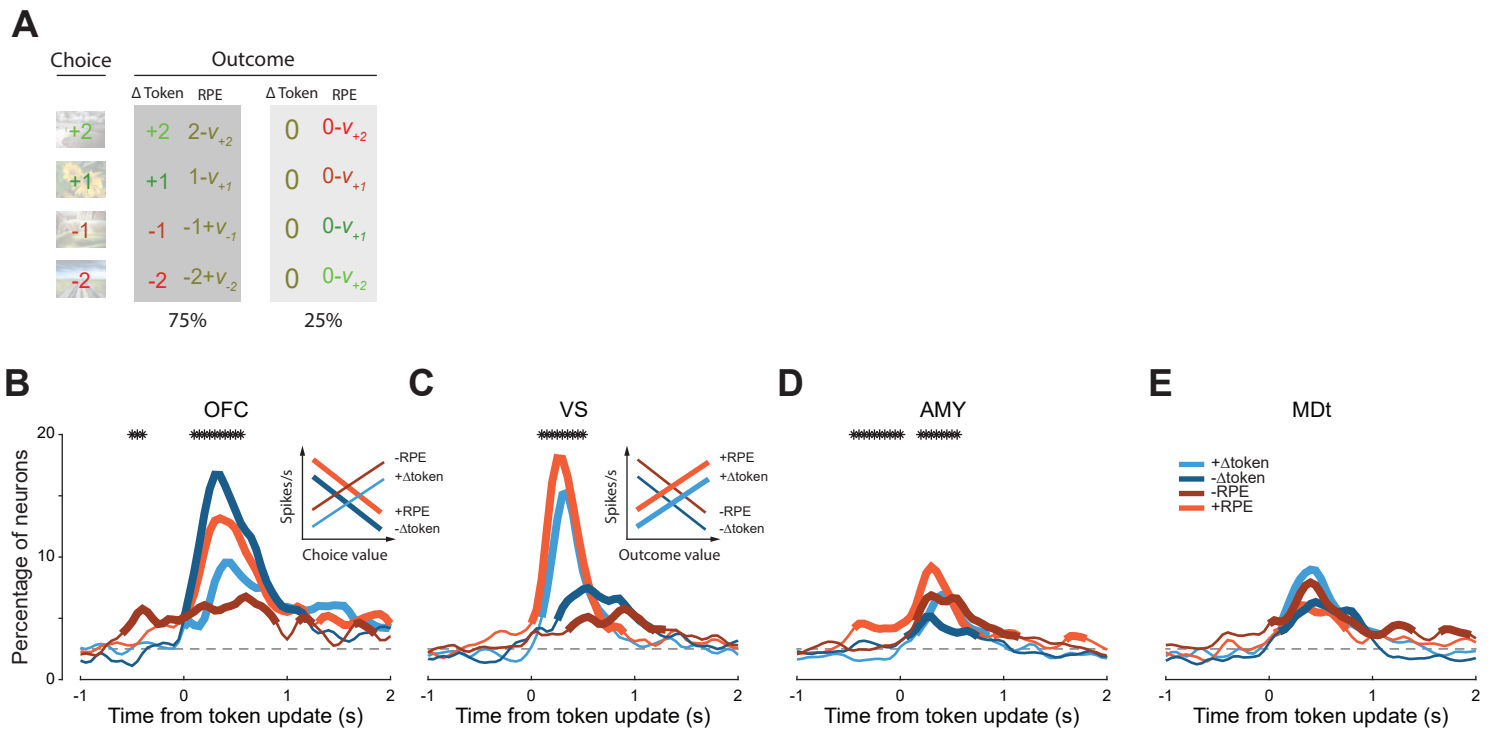


Figure 5

847 **Figure 5. Value updates in the ventral striatum and orbitofrontal cortex.**

848 (A) Relationship of *a priori* value, Δtoken , and RPE. v_h indicates the RW model-estimated value of the
849 image with *a priori* value h , which got closer to h as learning progressed.

850 (B-E) Percentage of neurons in each area encoding $+\Delta\text{token}$, $-\Delta\text{token}$, $+\text{RPE}$, and $-\text{RPE}$. Dashed horizontal
851 lines represent chance levels. Thick lines indicate a significant difference between the corresponding area
852 and chance level (binomial test, $p < 0.025$). The black asterisks indicate a significant difference among the
853 four categories (chi-square test, $p < 0.01$). Insets: OFC neurons encoded Δtoken and RPE in a manner
854 relevant to the value of choices. VS neurons encoded Δtoken and RPE correlate with the value of outcomes.
855 Thick lines represent more neurons.

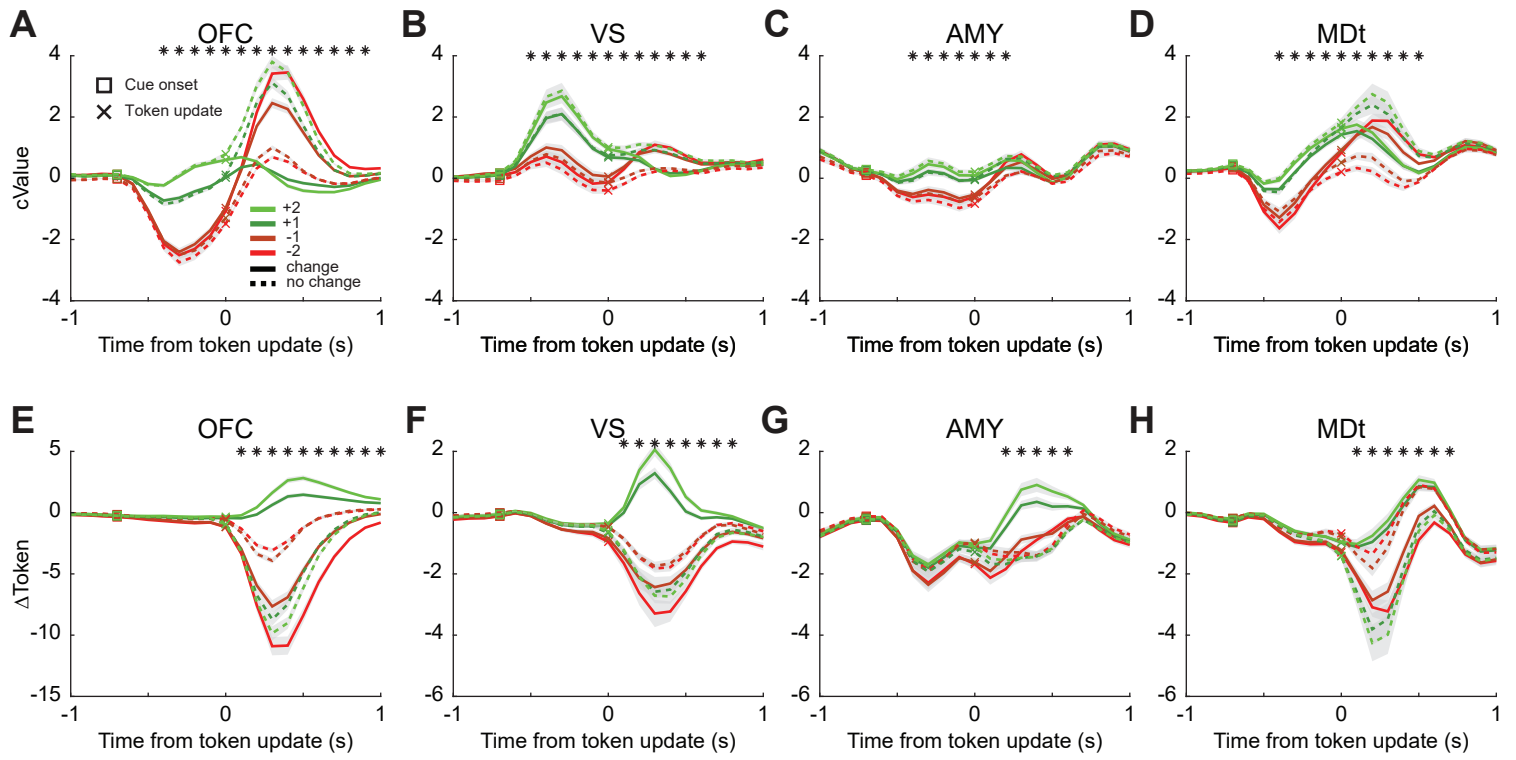


Figure 6

856 **Figure 6. Population representation of gains and losses**

857 (A-D) Projections of neural population activity on the cValue axis in the task-variable-specific one-
858 dimensional subspace. The y-axis indicates the value of the projection, divided by the number of neurons
859 (always 500). All neural activity was also z-scored. Therefore, the y-axis represents the average z-scored
860 population deviation per cell. Trials were grouped by cValue and outcome. Solid lines indicate the
861 conditions when tokens were delivered, and dashed lines indicate the conditions in which tokens were not
862 delivered. The analysis was repeated 1000 times by randomly choosing 500 neurons from each population
863 without replacement. Shaded zones show the mean \pm SD across iterations. The squares and crosses
864 represent the cue onset and token update. The black asterisks indicate a significant difference among the
865 eight conditions (bootstrap test, $p < 0.01$).

866 (E-H) Projections of neural population activity on the Δ token axis in the task-variable-specific one-
867 dimensional subspace.

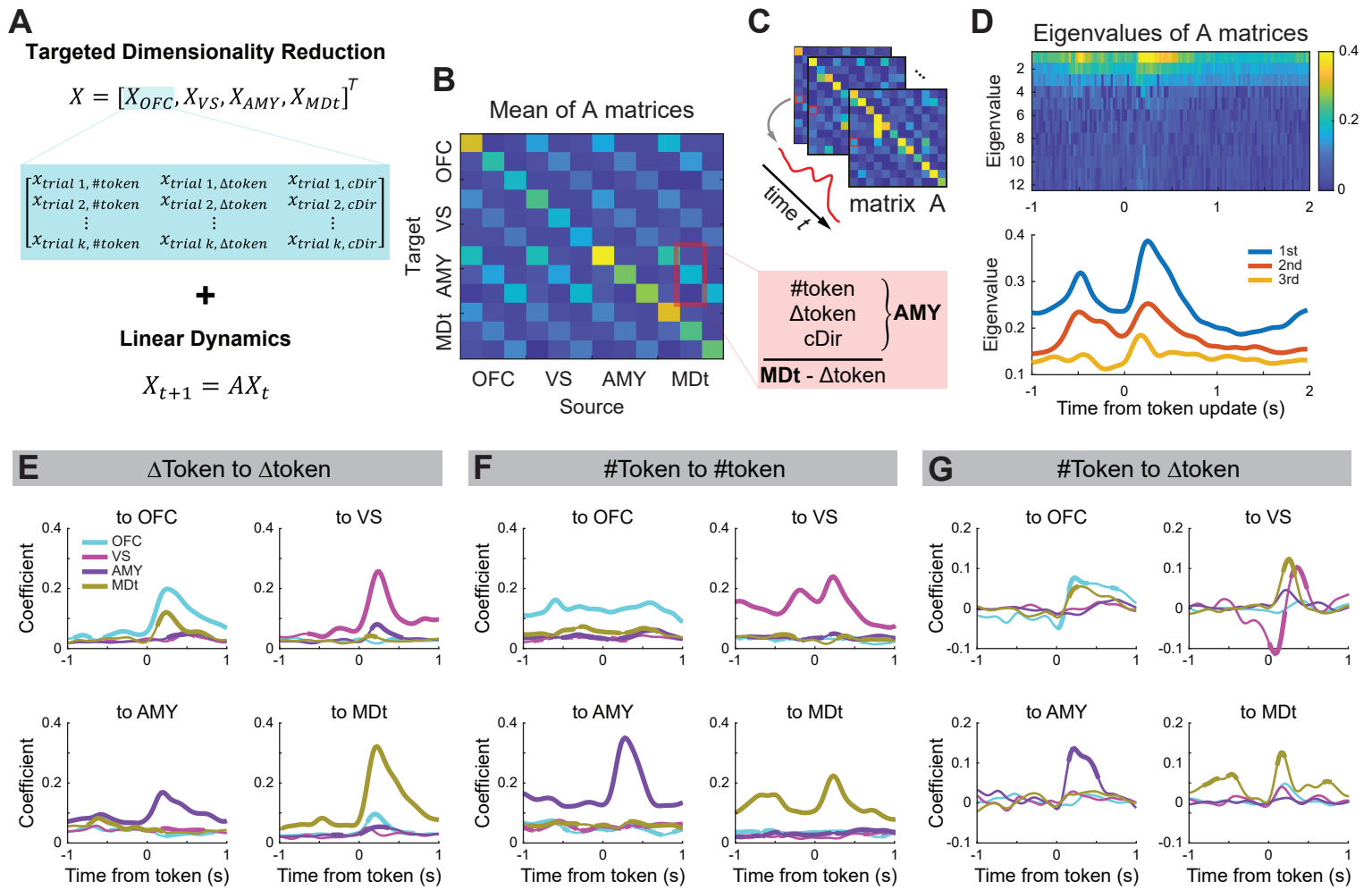


Figure 7

868 **Figure 7. Flow of token information in the ventral network.**

869 (A-D) Detecting network communications with low dimensional information dynamics. (A) The task-
870 related population response matrix consisted of single-trial-level population projections. Linear dynamic
871 regression measured time-dependent flows of information across task variables and brain areas. (B) The
872 loading matrix A maps the relationship between rows in matrix X . The x-axis represents the source areas
873 and variables, and the y-axis represents the target areas and variables. For example, the values indicated by
874 the red square represent the flow of Δ token information from the MDt to #token, Δ token, and cDir
875 information in the AMY. (C) The continuous values in A matrices represent the strength of information
876 flow across task variables and brain areas along the time course of the trial. (D) The eigenvalues of A
877 matrices capture the strength of the overall information flow among the areas in the network.

878 (E-G) The flow of token information across areas, including within the Δ token (E), #token (F), and between
879 them (G). The titles indicate the source and target task variables. Each panel indicates one target area, and
880 lines indicate source areas. Thick lines indicate a significant difference between the corresponding area and
881 shuffled data (two sides z-test, $p < 0.01$). Shuffling was done by keeping the same conditions but shuffling
882 the trial order of matrix X .