



# A quantitative reward prediction error signal in the ventral pallidum

David J. Ottenheimer<sup>1,6</sup>, Bilal A. Bari<sup>1,2,6</sup>, Elissa Sutlief<sup>1</sup>, Kurt M. Fraser<sup>3</sup>, Tabitha H. Kim<sup>3</sup>, Jocelyn M. Richard<sup>3,4</sup>, Jeremiah Y. Cohen<sup>1,2,5</sup> and Patricia H. Janak<sup>1,3,5</sup> ✉

**The nervous system is hypothesized to compute reward prediction errors (RPEs) to promote adaptive behavior. Correlates of RPEs have been observed in the midbrain dopamine system, but the extent to which RPE signals exist in other reward-processing regions is less well understood. In the present study, we quantified outcome history-based RPE signals in the ventral pallidum (VP), a basal ganglia region functionally linked to reward-seeking behavior. We trained rats to respond to reward-predicting cues, and we fit computational models to predict the firing rates of individual neurons at the time of reward delivery. We found that a subset of VP neurons encoded RPEs and did so more robustly than the nucleus accumbens, an input to the VP. VP RPEs predicted changes in task engagement, and optogenetic manipulation of the VP during reward delivery bidirectionally altered rats' subsequent reward-seeking behavior. Our data suggest a pivotal role for the VP in computing teaching signals that influence adaptive reward seeking.**

Adaptive behavior is characterized by responding flexibly to stimuli in our environments. The framework of reinforcement learning is a well-established approach for describing how individuals flexibly interact with their environments to maximize reward<sup>1</sup>. Reinforcement learning formalizes the notion that individuals integrate information about past rewards to make predictions about the future. Deviations from these predictions, known as RPEs, are used to iteratively update future predictions<sup>2</sup>. One remarkable extension of reinforcement learning to neuroscience was the discovery that midbrain dopamine neurons encode RPEs<sup>3</sup> and do so over local timescales<sup>4</sup>.

Although the midbrain dopamine system has been extensively studied for its role in learning, many limbic regions are involved in reward-guided behavior, and a role for RPE signals across this broader circuit has been largely overlooked. One candidate is the VP, a region critical for numerous reward-based behaviors<sup>5,6</sup>. The VP is part of the ventral striatopallidal system within the basal ganglia. Historically, the VP has been viewed primarily as an output of the nucleus accumbens (NAc), transmitting reward-related information from the NAc to downstream motor regions to mediate behavioral responses<sup>7</sup>. However, this view has recently been challenged, and the VP is now known to encode neural representations that are not directly inherited from the striatum<sup>8–11</sup>. This has renewed interest in the VP as an important site for reward processing in its own right. Large fractions of VP neurons are responsive to reward-related task events<sup>8–10,12,13</sup>, and there are hints that some VP neural activity is consistent with an RPE signal<sup>14–16</sup>. Moreover, phasic manipulations of VP activity impact reward-based behaviors<sup>8,15,17,18</sup>, suggesting links between the reward-related computations in the VP and behavioral responses.

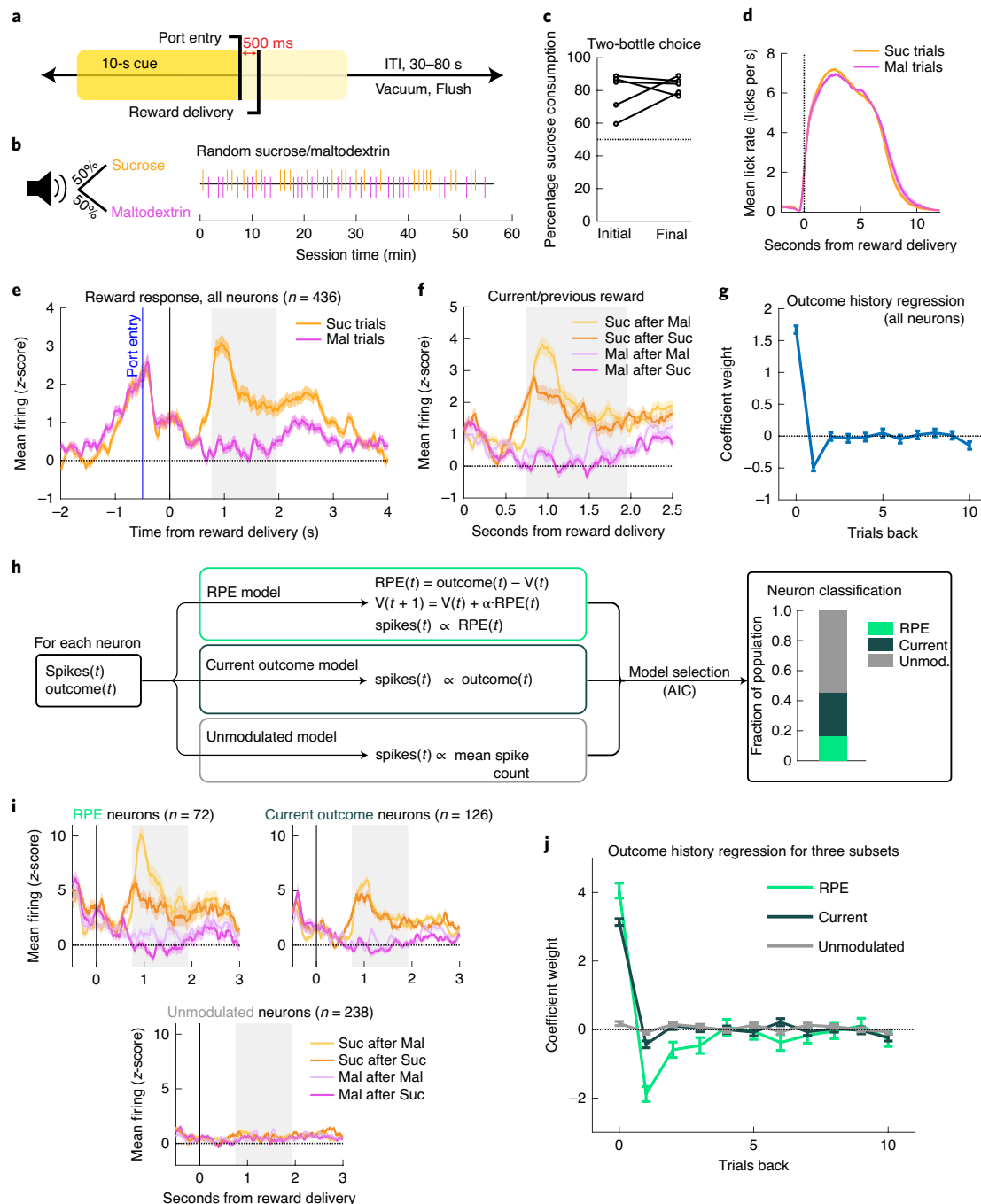
In the present study, by recording from the VP in rats performing a series of reward-seeking tasks, we demonstrate that VP neural activity is quantitatively consistent with an RPE signal. By adapting and fitting computational models to predict spike counts of individual

neurons, we classify a subset of VP neurons as RPE encoding. Importantly, we demonstrate that our RPE model predicts key features of VP neural activity, including RPE tuning and trial-by-trial firing rates, in contrast to poorer prediction of these features by the NAc, a main input to the VP. We further find that VP RPE neuron activity predicts subsequent task engagement, and optogenetic manipulation of the VP bidirectionally impacts task engagement.

## Results

**Preference-based reward prediction errors in the VP.** Rats were trained to respond to a 10-s white noise cue indicating the availability of 10% solutions of either sucrose or maltodextrin, contingent on entry into the reward port (Fig. 1a). In this task (random sucrose/maltodextrin), there was only one cue, which predicted sucrose or maltodextrin reward with equal probability. This task design ensured that rats could not accurately predict upcoming reward identity (Fig. 1b). As reported<sup>9,19</sup>, rats preferred sucrose when given free access to both sucrose and maltodextrin in their home cage, despite the rewards' equivalent caloric value (Fig. 1c). Nevertheless, they licked robustly for both during the task (Fig. 1d), reflecting the high palatability of both outcomes. This feature allowed us to control for a contribution of motor responses to reward-specific neural activity. We recorded the activity of 436 VP neurons while rats ( $n = 5$ ) performed the task (Extended Data Fig. 1) (some analyses of these recordings have been published previously<sup>9</sup>). Despite similar licking patterns, sucrose and maltodextrin evoked significantly different neural responses, with a higher mean firing rate when sampling sucrose (Wilcoxon's signed-rank test on all neurons' mean firing 0.75–1.95 s after sucrose or maltodextrin delivery,  $P = 8 \times 10^{-38}$ ), consistent with a preference for sucrose (Fig. 1e). Moreover, the previous outcome modulated the reward signal in a direction consistent with RPE coding (Fig. 1f). For example, receiving sucrose on the previous trial increased the expectation of future sucrose, leading to decreased firing when sucrose was delivered on

<sup>1</sup>Solomon H. Snyder Department of Neuroscience, Johns Hopkins University, Baltimore, MD, USA. <sup>2</sup>Brain Science Institute, Johns Hopkins University, Baltimore, MD, USA. <sup>3</sup>Department of Psychological and Brain Sciences, Johns Hopkins University, Baltimore, MD, USA. <sup>4</sup>Department of Neuroscience, University of Minnesota, Minneapolis, MN, USA. <sup>5</sup>Kavli Neuroscience Discovery Institute, Johns Hopkins University, Baltimore, MD, USA. <sup>6</sup>These authors contributed equally: David J. Ottenheimer, Bilal A. Bari. ✉e-mail: [patricia.janak@jhu.edu](mailto:patricia.janak@jhu.edu)



**Fig. 1 | A subset of VP neurons signal preference-based RPEs. a**, Task structure: entering the reward port during a 10-s cue-triggered reward delivery. **b**, There was a 50:50 probability of receiving sucrose or maltodextrin solutions, as seen in the example session. **c**, Percentage sucrose of total solution consumption in a two-bottle choice, before (Initial) and after (Final) recording ( $n = 5$  rats). **d**, Mean  $\pm$  s.e.m. of the lick rate relative to pump onset ( $n = 25$  sessions from 5 rats). **e**, Mean  $\pm$  s.e.m. of the activity of all recorded neurons on sucrose (Suc) and maltodextrin (Mal) trials ( $n = 436$  neurons from 5 rats). The gray rectangle indicates the window used for analysis in **g**, **h** and **j**, and all equivalent analyses in subsequent figures. **f**, Mean  $\pm$  s.e.m. of the activity of all recorded neurons on trials sorted by previous and current outcome. **a, c–f**, Adapted from ref. <sup>9</sup>. **g**, Coefficients  $\pm$  s.e.m. from a linear regression fit to the z-scored activity of all neurons ( $n = 436$  neurons) and the outcomes on the current trial and preceding 10 trials. **h**, Schematic of model-fitting and neuron classification process. For each neuron, the reward outcome and spike count after reward delivery on each trial were used to fit three models: RPE, current outcome and unmodulated (unmod.). The AIC was used to select which model best fit each neuron's activity (right). **i**, Mean  $\pm$  s.e.m. of the activity of neurons best fit by each of the three models, plotted according to the previous and current outcome. **j**, Coefficients  $\pm$  s.e.m. for the outcome history linear regression for each class of neurons ( $n = 72$  RPE, 126 current outcome and 238 unmodulated neurons).

the current trial. The expected trend held true for all combinations of past/current outcomes, suggesting that VP neural activity may contain an RPE signal.

Intrigued by the possibility of RPE signaling in the VP, we expanded on these initial findings by quantifying the impact of current and previous outcomes on reward-evoked firing in the VP.

We applied a linear regression that has previously been used to quantify the effect of reward history on dopamine neuron firing<sup>4</sup>. When the activity of all neurons was pooled, only the current trial and previous trial impacted firing rates at the time of the outcome (Fig. 1g). Although this pattern is consistent with RPE coding, it is on a much shorter timescale than has been observed for dopamine neurons<sup>4,20</sup>. One limitation of the pooled linear regression approach is that it assumed that the VP is largely homogeneous, which risks introducing bias into coefficient estimates. This left open the possibility that the VP contains subsets of neurons encoding reward history on longer timescales, and led us to analyze individual neuronal responses.

To identify neurons in the VP sensitive to reward history, we developed three computational models to fit firing rates of individual neurons, corresponding to three potential patterns of neuronal activity. The first model, 'RPE', is based on the Rescorla–Wagner model<sup>2</sup>. The model generated trial-by-trial value estimates ( $V$ ) which constitute reward predictions. On each trial, an RPE was generated by the difference between actual and predicted rewards, and this RPE was multiplied by a learning rate ( $\alpha$ ) before updating  $V$  for the next trial. Small learning rate values allow for integration of the reward history multiple trials into the past. The spike counts of individual neurons were fit to the estimated RPE on each trial (Fig. 1h). We also fit two additional models to serve as controls, one in which the spike count was determined only by the current outcome and one with no impact of outcome (unmodulated). We used maximum likelihood estimation to fit the models to each neuron and selected the most parsimonious model using the Akaike information criterion (AIC), which selects the best-fit model after penalizing for model complexity. This classification process revealed that 17% of VP neurons were best described by the RPE model, and another 29% were best fit by the current outcome only (Fig. 1h); notably, of the 47% of neurons we had previously classified as sucrose preferring in our previous work<sup>9</sup>, 74% were classified as either RPE or current outcome in the present study, demonstrating general agreement between the approaches.

We plotted the mean reward-evoked activity of each subset of neurons according to previous and current outcome, and found that the firing rates of each subset matched the properties of the models with which they were classified (Fig. 1i). For instance, 86% of RPE neurons had both higher firing rates for sucrose after maltodextrin than for sucrose after sucrose (positive RPE), and lower firing rates for maltodextrin after sucrose than for maltodextrin after maltodextrin (negative RPE). We then performed the same linear regression of outcome history on each subset of neurons rather than the entire population; this revealed an exponential decay-like influence of multiple previous trials on firing of neurons best fit by the RPE model, indicating that VP neurons modulated by reward history were in fact integrating information over a more extended period of time (Fig. 1j, non-zero weights for one to three trials back). Indeed, the mean (median) learning rate across all neurons was 0.56 (0.52); this corresponds to an exponential learning process with a half-life of 0.84 (0.94) trials, indicating that neurons accumulate information over  $\sim 4.22$  (4.72) trials to reach a steady-state value estimate (full distribution given in Extended Data Fig. 2). Thus, given the closely matched caloric value and motor responses to each reward, these data indicate that some VP neurons signal an outcome history-based RPE according to reward preference.

**The VP encodes RPEs more robustly than the NAc.** We next asked how faithfully VP neurons encoded RPEs. Our fitting procedure allowed recovery of trial-by-trial estimates of RPEs, based on parameter estimates for that individual neuron, as well as the outcome history for that session. Model-derived RPEs were strongly correlated with the activity of individual neurons (Fig. 2a) and the average activity across all RPE neurons (Fig. 2b). Importantly, this approach revealed a finer dynamic range of firing than revealed by only

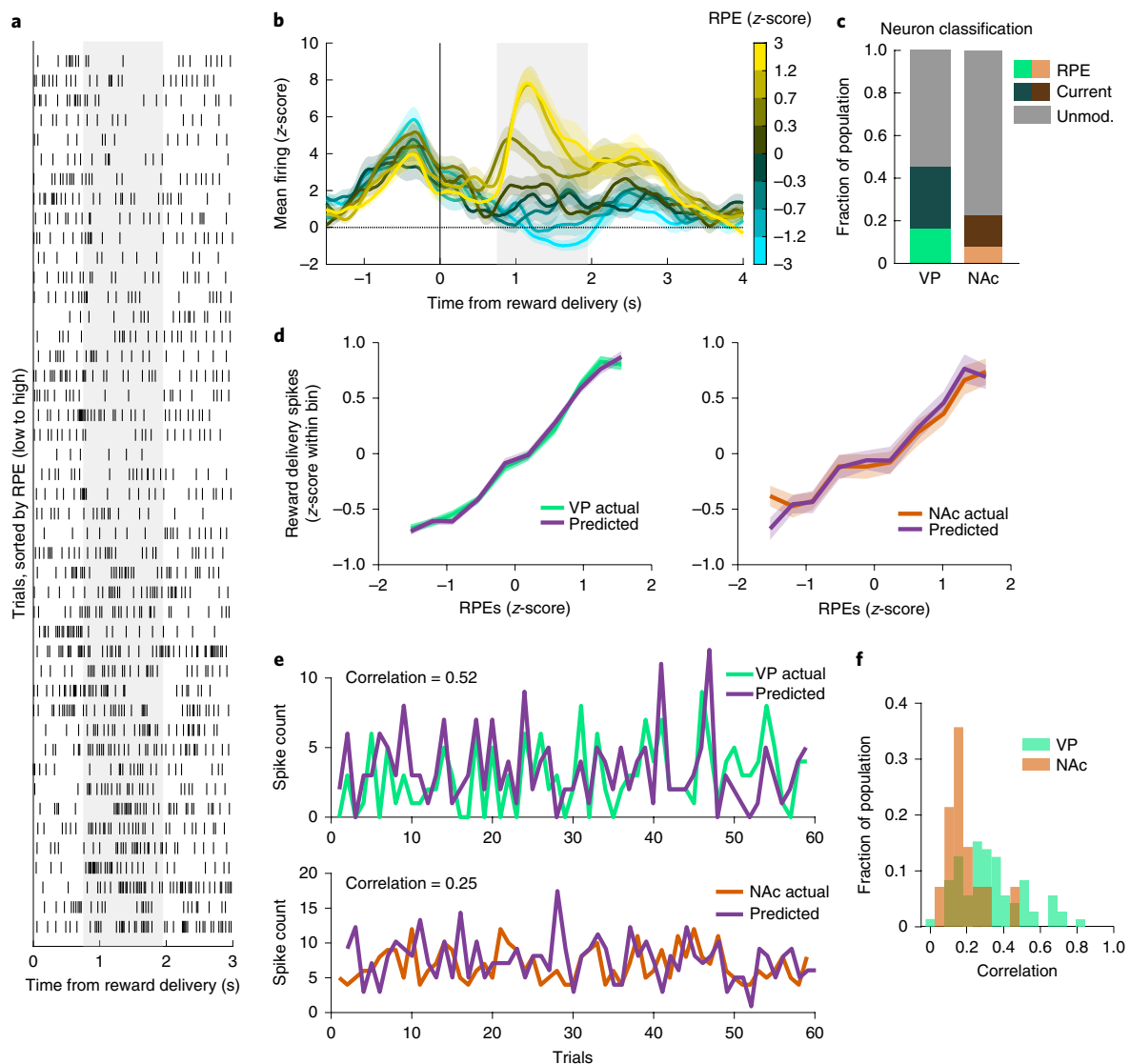
looking at current and previous outcomes (Fig. 1i). We next generated RPE tuning curves for these neurons and found a strong monotonic relationship ( $t_{3,961} = 41.3$ ,  $P = 2 \times 10^{-310}$ , linear relationship between RPEs and  $z$ -scored firing rates). As a stronger test, we used parameters estimated for each neuron to simulate RPE-correlated spike counts and generated an 'ideal' RPE tuning curve. We observed a clear overlap between real and simulated tuning curves (Fig. 2d). Finally, we quantified the correlation between predicted and real spike counts and found good agreement (Pearson's correlation coefficient: mean 0.34, median 0.31; Fig. 2e,f).

To contextualize the robustness of VP RPE responses, we ran the same analysis on neurons ( $n = 183$  neurons from 6 rats) recorded during the same task in the NAc<sup>9</sup> (Extended Data Fig. 1), a major VP input<sup>6,7</sup>. We found fewer cells with activity best fit by the RPE model in the NAc than in the VP and by the current outcome model (Fig. 2c). Moreover, NAc neurons classified as RPE signaling were described less well by the model than similarly classified VP neurons. This was evident by a poorer match between real and simulated neuron tuning curves (mean squared error between real and simulated tuning curves; bootstrapped 95% confidence intervals (CIs): (1.24 1.38) in the VP, (1.44 1.77) in the NAc; Fig. 2d) and in poorer correlation between model-predicted and actual spiking for individual RPE neurons (Pearson's correlation coefficient: mean 0.18, median 0.15; Wilcoxon's rank-sum test,  $P = 0.0007$ ; Fig. 2e,f). Thus, in the current task, the VP has more robust RPE signaling than the NAc, a notable finding because the VP is purported to inherit its firing from the NAc.

### An expanded value space reveals stronger RPE signaling in VP.

One shortcoming of the experiment contrasting sucrose and maltodextrin is that the similar palatability of the outcomes may not fully probe the limits of value signaling, potentially constraining our ability to identify RPE neurons; maltodextrin delivery does not typically strongly inhibit responses at the time of reward (Fig. 1e). We previously found that delivering water, an outcome less rewarding than maltodextrin, more strongly inhibited firing rates ('random sucrose/maltodextrin/water' task; Fig. 3a–c)<sup>9</sup>. We hypothesized that this expansion of the dynamic range of firing would reveal additional RPE neurons. We applied the same models as before to neurons recorded during this three-outcome task ( $n = 254$ ) to identify cells with firing patterns that reflected outcome history-based RPEs, current outcome only or no modulation, with an additional free parameter to estimate the value associated with maltodextrin on the scale of water (0) to sucrose (1). As hypothesized, a greater proportion of neurons was best fit by the RPE model in this task than in the random sucrose/maltodextrin task (29% versus 17%,  $\chi^2 = 15.3$ ,  $P = 9 \times 10^{-5}$ ; Fig. 3d). Outcome history regressions revealed an impact of many previous trials on these neurons (Fig. 3e, non-zero weights for one to three, five to seven, and nine trials back). We observed graded changes in firing rates as a function of estimated RPEs for individual neurons (Fig. 3f); this relationship was consistent in the population-average peristimulus time histogram (PSTH) (Fig. 3g). Firing rates of these RPE neurons monotonically increased as a function of estimated RPEs, and this relationship was consistent with tuning curves for simulated RPE neurons (Fig. 3h). Moreover, the model's predictions of trial-by-trial spiking for each neuron was robust and stronger than in the random sucrose/maltodextrin task (Pearson's correlation coefficient: mean 0.48, median 0.49; Wilcoxon's rank-sum test between the VP–RPE correlation in the 'random sucrose/maltodextrin' and that in the 'random sucrose/maltodextrin/water' task,  $P = 3 \times 10^{-5}$ ; Fig. 2e–f). Thus, with outcomes spanning an expanded value space, we found more neurons that encode RPEs and do so more robustly.

**VP reward activity mediates trial-by-trial task engagement.** The presence of adaptive RPE signals in the VP in these tasks raises the



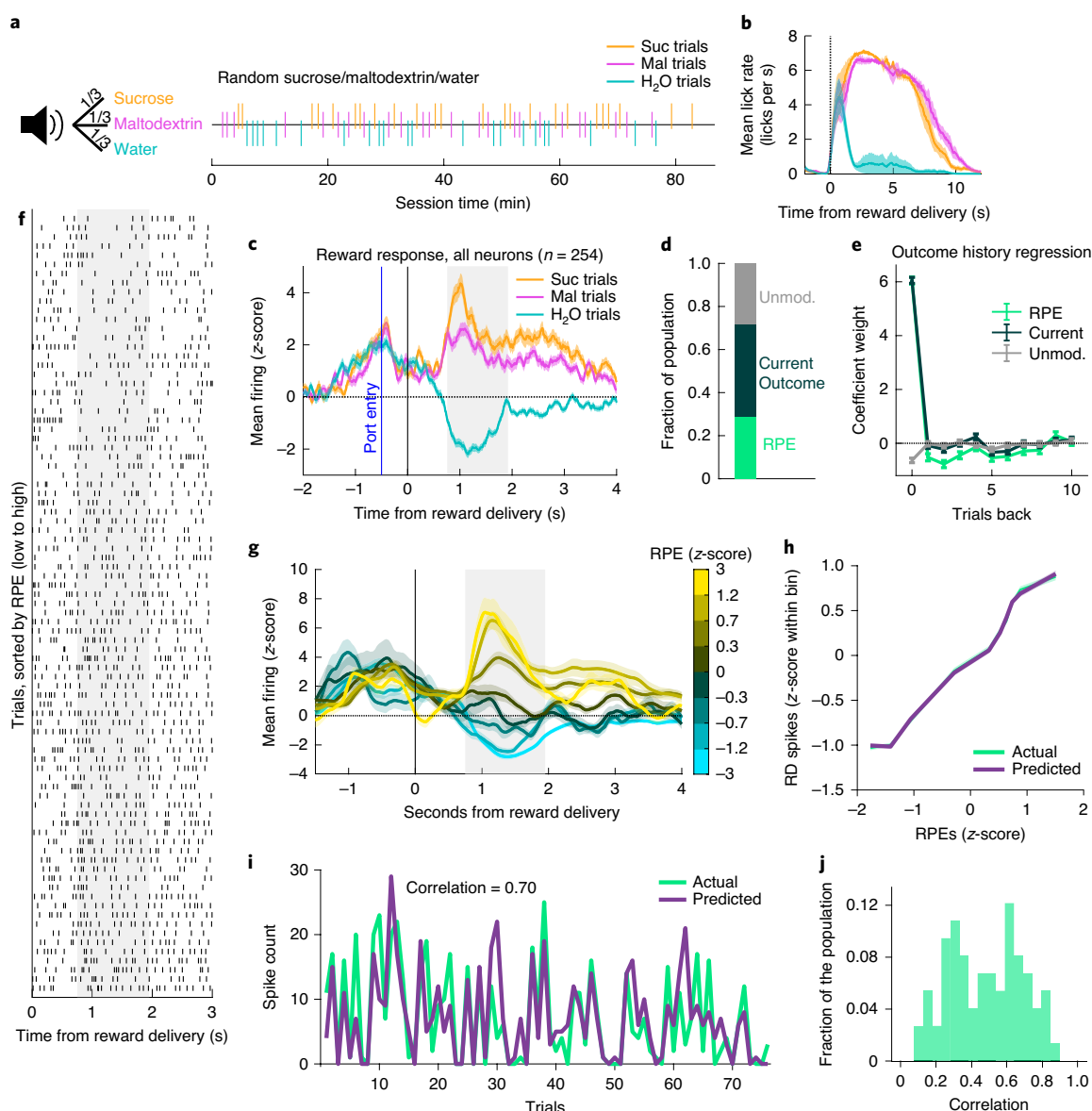
**Fig. 2 | RPE encoding is more prevalent and robust in the VP than in the NAc.** **a**, Raster of an individual VP neuron's spikes on each trial, aligned to reward delivery, and sorted by the model-derived RPE value for each trial. The gray-shaded region indicates window used for analysis. **b**, Population mean  $\pm$  s.e.m. of all VP RPE neurons identified in Fig. 1. The trials for each neuron are binned according to their model-derived RPE. **c**, Proportion of the population in the VP and NAc classified as RPE, current outcome or unmodulated. There were fewer RPE cells in the NAc than in the VP (8% versus 17%,  $\chi^2 = 8.3$ ,  $P = 0.004$ ) and current outcome cells (14% in the NAc versus 29% in the VP,  $\chi^2 = 13.6$ ,  $P = 0.0002$ ). **d**, Mean  $\pm$  s.e.m. of population activity of simulated and actual RPE neurons according to each trial's RPE value for the VP (top) and the NAc (bottom). **e**, The model-predicted and actual spike counts on each trial for one RPE neuron each from the VP (left) and the NAc (right). These neurons were the 85th percentile for correlation for each respective region. **f**, Distribution of correlations between model-predicted and actual spiking for all RPE neurons from each region.

question of whether rats were similarly adapting their behavior in response to reward outcomes. As the rats were freely moving, task participation represented a trade-off between reward seeking and competing interests, including rest, grooming and exploring the behavioral chamber. To evaluate task participation, we analyzed videos ( $n = 4$ ) from the random sucrose/maltodextrin/water recording sessions. To estimate trial-by-trial task engagement, we calculated the average distance from the port in each intertrial interval (ITI). This analysis revealed instances where rats traveled far from the reward port and, in some cases, remained far from the reward port at the start of the next trial (Fig. 4a). Consistent with their relative reward palatability, rats remained close to the reward port during the ITI after sucrose, moved further from the port after maltodextrin and even further after water delivery (Fig. 4b).

Next, we determined whether this behavioral measure, proximity to the reward port, was related to the VP neural signals we

characterized. Consistent with the idea that VP reward signals guide task engagement, there was on average a negative correlation between the activity of VP RPE cells ( $n = 74$ ) and current outcome cells ( $n = 108$ ) at reward delivery and distance from the port during the next ITI (Fig. 4c,d); there was a modest positive correlation for unmodulated cells, perhaps reflecting a VP population that promotes avoidance<sup>15,17,18</sup>. The negative correlation for RPE and current outcome cells indicates that rats traveled around the chamber and remained far from the reward port after the activity of these neurons was low (that is, negative prediction errors or less-preferred outcomes); conversely, rats remained closer to the port after high activity. This pattern of results held for sessions contrasting sucrose and maltodextrin as well: rats were closer to the port after sucrose trials ( $P = 0.048$ , Wilcoxon's signed-rank test), and there was a negative correlation between RPE and current outcome cell activity and distance from the port ( $P = 0.001$  for both, Wilcoxon's signed-rank



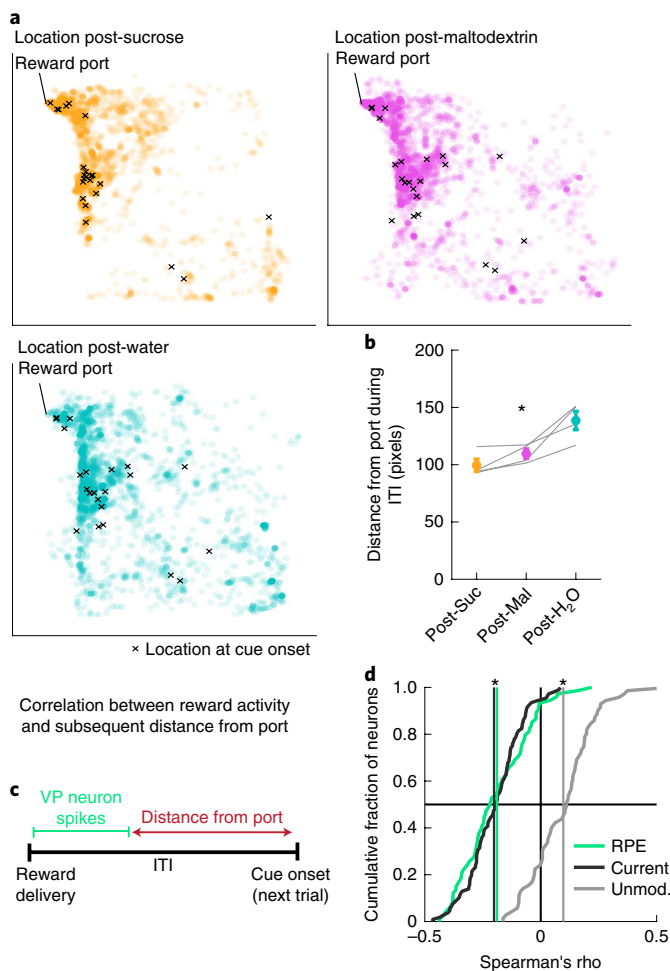


**Fig. 3 | An expanded value space reveals stronger RPE signaling in the VP.** **a**, Task structure: on each trial, there was a one in three probability each of receiving sucrose, maltodextrin or water, as seen in the example session. **b**, Mean  $\pm$  s.e.m. of lick rate relative to pump onset ( $n=4$  sessions from 3 rats). **c**, Mean  $\pm$  s.e.m. of activity of all recorded neurons on sucrose, maltodextrin and water trials ( $n=254$  neurons from 3 rats). **d**, Fraction of the population of neurons recorded in this task best fit by each of the three models. **e**, Coefficients  $\pm$  s.e.m. for outcome history regression for each of the three classes of neurons ( $n=74$  RPE, 108 current outcome and 72 unmodulated neurons). **f**, Raster of an individual neuron's spikes on each trial, aligned to reward delivery and sorted by the model-derived RPE value for each trial. The gray-shaded region indicates the window used for analysis. **g**, Population mean  $\pm$  s.e.m. of all RPE neurons. The trials for each neuron are binned according to their model-derived RPE. **h**, Mean  $\pm$  s.e.m. of the population activity of simulated and actual VP RPE neurons according to each trial's RPE value. **i**, The model-predicted and actual spike counts on each trial for the RPE neuron with the 85th percentile correlation. **j**, Distribution of correlations between model-predicted and actual spiking for all RPE neurons.

test), but not for unmodulated cells ( $P=0.75$ , Wilcoxon's signed-rank test).

The correlation between VP neuron activity and task engagement suggests that the VP may causally influence task engagement. To explore this possibility, we used an optogenetic approach. Rats were infused with virus containing either the inhibitory opsin, ArchT3.0-eYFP ( $n=7$ ), or enhanced yellow fluorescent protein (eYFP) alone as a control ( $n=7$ ), and implanted with optic fibers aimed at the VP (Fig. 5a and Extended Data Fig. 3). We then trained these rats on a simplified task in which port entry during a 10-s cue earned a sucrose reward. On half the trials, we inhibited the VP for 5 s starting at onset of sucrose delivery, mimicking a negative

prediction error or a nonpreferred outcome signal (Fig. 5b). Similar to water delivery (a less-preferred option), optogenetic inhibition of the VP increased rats' average distance from the port during the next ITI relative to control rats ( $P=0.01$ , Wilcoxon's rank-sum test) (Fig. 5c–e). We then performed the complementary experiment by injecting channelrhodopsin (ChR2)-containing virus ( $n=10$ ) or green fluorescent protein (GFP) control ( $n=7$ ) into another group of rats (Fig. 5f and Extended Data Fig. 3). Rats were trained on the same task, and on half the trials we stimulated the VP for 2 s at 40 Hz, approximating a positive prediction error or a preferred outcome (Fig. 5g). VP stimulation increased subsequent task engagement, decreasing the distance during the next ITI relative to



**Fig. 4 | VP reward activity tracks changes in trial-by-trial task engagement.** **a**, All locations of a rat from an example session during the ITI after sucrose (left), maltodextrin (right) and water (bottom) delivery. Each circle is one location during a 0.2-s bin. X marks the location at cue onset for the subsequent trial. The chamber is 32.4 × 32.4 cm<sup>2</sup> (approximately 306 × 306 pixels<sup>2</sup>). **b**, Mean ± s.e.m. of the distance from the port during the ITI after sucrose (orange), maltodextrin (pink) and water (blue) trials during recording sessions ( $n = 4$  sessions from 3 rats). The gray lines represent the mean for one subject in one session ( $*r = -0.86$ ,  $P = 0.0004$ , Spearman's rank correlation coefficient between distance from port and reward preference ranking). **c**, Approach for correlating the activity of individual VP cells with distance from the port on a trial-by-trial basis. **d**, Distribution of correlations between individual VP neurons' firing rates on each trial and the distance from the port during the subsequent ITI (\*significant shift in mean correlation coefficient (vertical lines) compared with 1,000 shuffles of data for RPEs ( $P = 8 \times 10^{-10}$ ), current outcome ( $P = 3 \times 10^{-18}$ ) and unmodulated neurons ( $P = 0.00008$ ) (Wilcoxon's signed-rank test, two sided)).

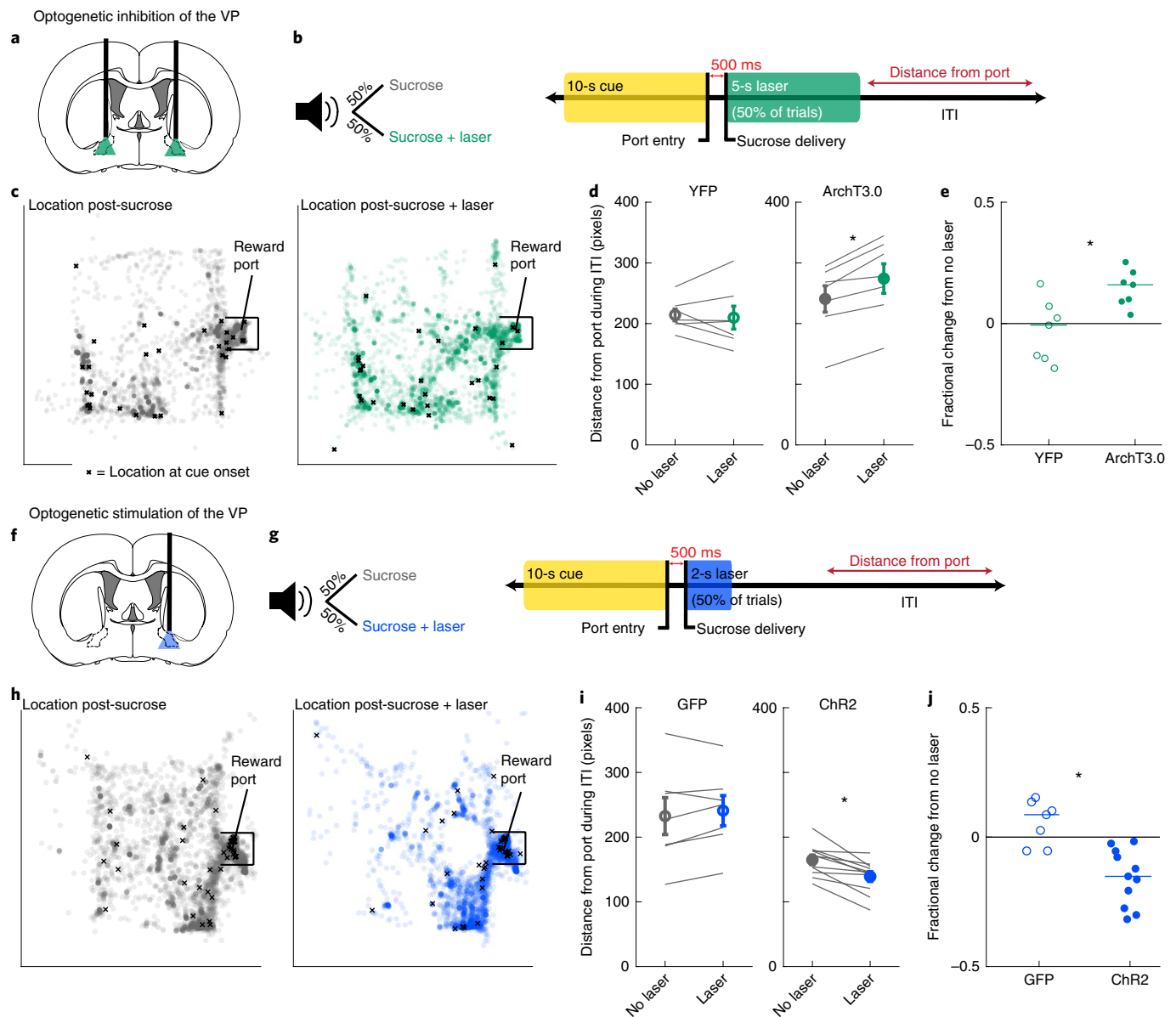
control rats ( $P = 0.001$ , Wilcoxon's rank-sum test) (Fig. 5h–j). Thus, VP activity is instructive for task engagement-related behavior, suggesting that outcome-related signals in the VP are used to motivate task performance.

To ensure the robustness of these findings, we more closely analyzed the effect of optogenetic manipulation on behavior. First, we analyzed how the laser affected ongoing behavior. There was no noticeable impact of the laser on time spent in the reward port for the control groups nor for the ArchT3.0 group, meaning that the effects of manipulation on task engagement were probably not due to interruption of the consumption phase (Extended Data Fig. 4a).

It is interesting that, in the Chr2 group, stimulation of the VP caused the rats to move their head out of the port, leading to delayed reward consumption (Extended Data Fig. 4b). To control for the possibility that delayed consumption led to the difference in ITI distance, we restricted our analysis to the ITI period beyond 15 s after reward delivery, when port entry occupancy on laser and no-laser trials was similar ( $P = 0.12$ , Wilcoxon's signed-rank test). We confirmed that Chr2 rats were closer to the port after laser trials than after no laser, relative to the control ( $P = 0.01$ , Wilcoxon's rank-sum test). Finally, we ran a similar experiment to find out how VP activation during the cue (for 2 s), rather than reward, would influence behavior (Extended Data Fig. 4c) and observed no effect on ITI distance in Chr2 rats relative to the control ( $P = 0.36$ ; Extended Data Fig. 4d,e). This indicates that VP activity during the reward outcome epoch specifically influences task engagement during the subsequent ITI.

**VP RPE neuron firing adapts to repeated reward presentations.** Repeated presentation of the same reward (or sets of rewards) can produce adaptation in neural responses as the outcome becomes expected<sup>21–23</sup>, a phenomenon that can be explained by RPE models. We investigated whether VP neurons also attenuate their reward-evoked firing to repeated outcomes by analyzing the activity of neurons ( $n = 348$ ) recorded during a variation of the sucrose and maltodextrin task in which each reward was presented in blocks of 30 trials (Fig. 6a,b). Neural activity was fit to the same three models (RPE, current outcome and unmodulated; Fig. 1h), revealing a similar fraction of RPE neurons during this task as in the random sucrose/maltodextrin task (Fig. 6c). There were noticeable differences in the average firing rate of RPE neurons in the blocked task compared with the interspersed task, consistent with an acquired reward expectation in the blocked task (Fig. 6d,e). To determine how the reward-evoked activity evolved across each block, we plotted the activity in three-trial bins evenly spaced throughout the session (Fig. 6f–h). RPE neurons demonstrated notable reward-specific adaptations: a reduction in activity within sucrose blocks ( $t_{1,804} = -5.7$ ,  $P = 1 \times 10^{-8}$  for a linear model fitting RPE neuron activity to session progress in sessions with sucrose block first;  $t_{1,882} = -8.5$ ,  $P = 5 \times 10^{-17}$  for sucrose block second) and an increase within the maltodextrin block when maltodextrin was second ( $t_{1,697} = 4.3$ ,  $P = 0.00002$ ), although not when it was first ( $t_{1,821} = 0.38$ ,  $P = 0.71$ ), resulting in a significant interaction between the effects of session progress and outcome on the firing rates of RPE neurons (sucrose first:  $t_{1,1501} = -6.8$ ,  $P = 2 \times 10^{-11}$ ; sucrose second:  $t_{1,1703} = -6.4$ ,  $P = 2 \times 10^{-10}$ ). This pattern was consistent with predictions of the RPE model (Fig. 6f); over time, positive prediction errors for sucrose decrease as the expectation for sucrose builds across repeated sucrose trials and, similarly, negative prediction errors for maltodextrin attenuate (resulting in increased firing rate) as the expectation for maltodextrin builds across the maltodextrin block. Notably, current outcome and unmodulated neurons in the blocked sessions did not follow this pattern (Fig. 6f–h; all  $P > 0.05$  for interaction between session progress and outcome). RPE neurons from random sucrose/maltodextrin sessions also did not follow this pattern (Fig. 6f;  $t_{1,3959} = 1.7$ ,  $P = 0.08$ ), demonstrating that these across-session changes are specific to the blocked structure. Therefore the same RPE model that describes neurons sensitive to outcome history when rewards are randomly interspersed can also identify neurons in the VP that exhibit adaptation across blocks.

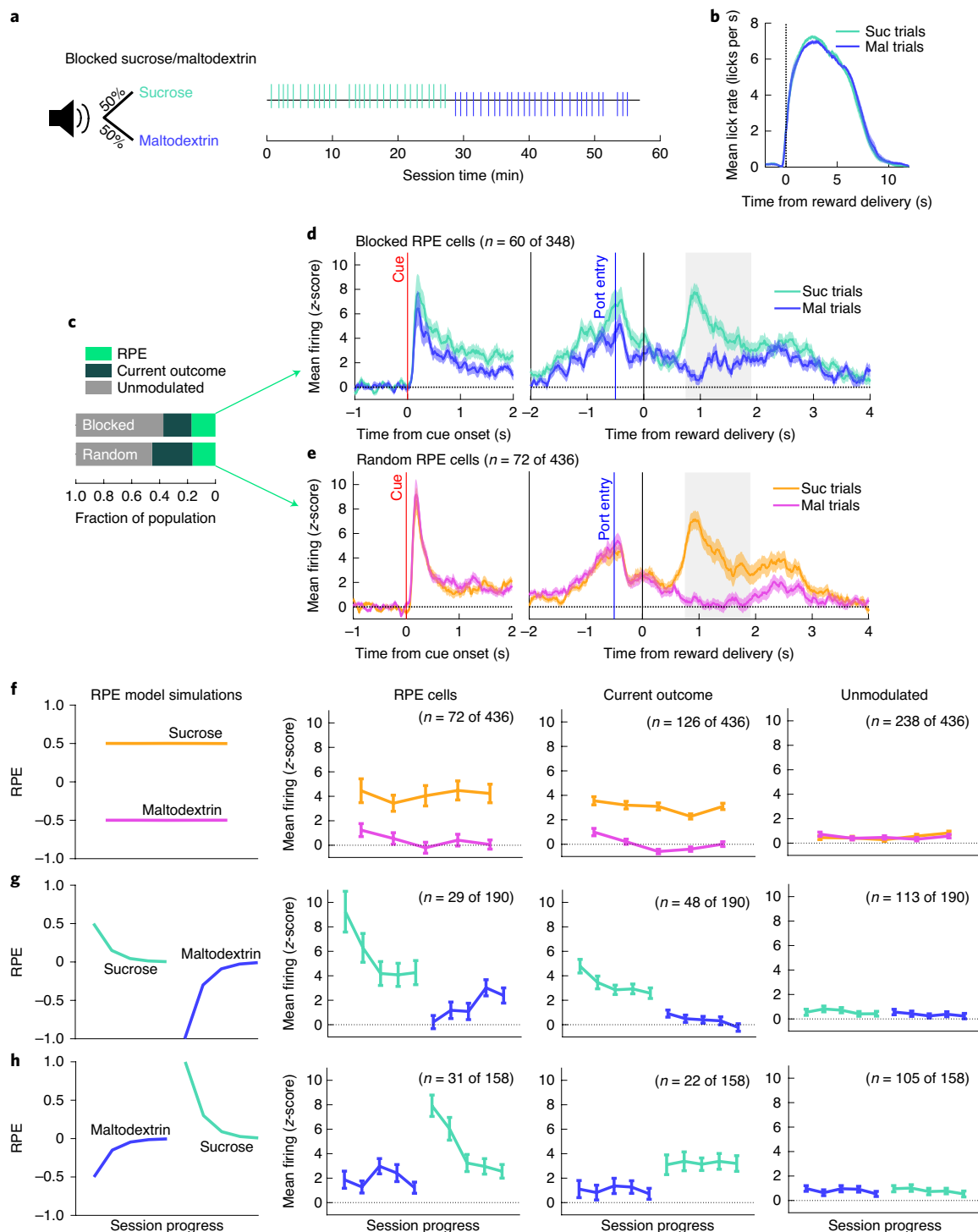
**Impact of reward-predicting cues on VP firing.** In the analysis thus far we applied a Rescorla–Wagner trial-based RPE model<sup>2,4</sup> to characterize how expectation is updated iteratively by the outcome of each trial. A critical expansion of the Rescorla–Wagner model is the temporal difference (TD) model, which allows within-trial updating of expectation by events such as reward-predicting cues<sup>1,3,24,25</sup>. We used a number of approaches to evaluate whether reward-predicting



**Fig. 5 | Manipulation of VP reward activity bidirectionally alters task engagement.** **a**, Optogenetic inhibition of VP with ArchT3.0. **b**, Experimental approach to evaluate the contribution of the VP to task engagement. On 50% of the trials, rats received green laser-mediated inhibition for 5 s. **c**, All locations of a rat from an example session during the ITI after sucrose delivery without laser (left) and with laser (right). Each circle is one location during a 0.2-s bin. X marks the location at cue onset for the subsequent trial. The chamber is 29.2 × 24.4 cm<sup>2</sup> (approximately 542 × 460 pixels<sup>2</sup>). **d**, Mean ± s.e.m. of the distance from the port in the ITI after sucrose with and without laser for animals receiving a control virus (YFP, left, *n* = 7 rats) or the ArchT3.0 virus (right, *n* = 7 rats). Individual rats' data are shown by the gray lines (\**P* = 0.02, Wilcoxon's signed-rank test, two sided). **e**, Fractional change in ITI distance from port for each rat (median: −0.01 YFP, *n* = 7 rats; 0.15 ArchT3.0, *n* = 7 rats; \**P* = 0.01, Wilcoxon's rank-sum test, two sided). **f**, Optogenetic stimulation of the VP with ChR2. **g**, As in **b**, but with 2 s of blue laser-mediated stimulation at 40 Hz, 10-ms pulse width. **h**, Example session, as in **c**. **i**, Mean ± s.e.m. of the distance from the port in the ITI after sucrose with and without laser for animals receiving a control virus (GFP, left, *n* = 7 rats) or the ChR2 virus (right, *n* = 11 rats). Individual rats' data are shown by the gray lines (\**P* = 0.001, Wilcoxon's signed-rank test, two sided). **j**, Fractional change in ITI distance from port for each rat (median: 0.09 GFP, *n* = 7 rats; −0.14 ChR2, *n* = 11 rats; \**P* = 0.001, Wilcoxon's rank-sum test, two sided).

cues impacted VP firing in a TD-like pattern. First, we analyzed the cue-evoked activity in the random sucrose/maltodextrin and random sucrose/maltodextrin/water tasks. Although there is only one cue in these sessions, the cue value can change according to the recently received outcomes. We used the same model-fitting classification procedure that we applied to firing at reward delivery to characterize cue-evoked activity. We compared the fits of the unmodulated model and a 'value' model, which is identical to the RPE model but maps *V* (the expected value) rather than RPE on

to the neuron's cue-evoked spiking on each trial (Extended Data Fig. 5). Although there were few neurons classified as encoding value in the sucrose/maltodextrin task, they were more common among neurons encoding RPEs at the outcome (15%) than non-RPE neurons (8%,  $\chi^2 = 4.2$ , *P* = 0.04). In the sucrose/maltodextrin/water task, in which a greater fraction of VP neurons encoded RPE at the outcome, there was also a greater fraction encoding value at the cue ( $\chi^2 = 5.9$ , *P* = 0.016; Extended Data Fig. 6) and, again, RPE neurons were more likely to encode value during the cue than non-RPE



**Fig. 6 | VP RPE neuron signaling adapts across reward blocks.** **a**, Task structure: sucrose and maltodextrin were presented in blocks of 30 trials, as seen in the example session. **b**, Mean  $\pm$  s.e.m. of the lick rate relative to pump onset ( $n = 14$  sessions from 5 rats). **c**, Proportion of neurons best fit by each of the three models in the random and blocked sucrose-maltodextrin tasks. **d**, Mean  $\pm$  s.e.m. of the activity of all RPE neurons from the blocked tasks aligned to cue onset and reward delivery ( $n = 60$  neurons from 5 rats). **e**, Mean  $\pm$  s.e.m. of the activity of all RPE neurons from the random sucrose-maltodextrin task aligned to cue onset and reward delivery ( $n = 72$  neurons from 5 rats). **f**, RPE model simulations (left) and mean  $\pm$  s.e.m. of the activity of RPE, current outcome and unmodulated cells from the random sucrose-maltodextrin task, plotted in bins of three trials evenly spaced throughout all completed sucrose and maltodextrin trials. **g**, As in **f**, for blocked sessions with sucrose first. **h**, As in **f** and **g** for blocked sessions with maltodextrin first.

neurons (value coding in 28% of RPE neurons, 9% of non-RPE neurons;  $\chi^2 = 14.8$ ,  $P = 0.0001$ ). This pattern of results indicates that some VP cells fire in a TD-like pattern, with the expected outcome reflected in the firing rate of both cue- and outcome-evoked activities.

A key demonstration of TD RPEs in dopamine neurons was the observation that firing at both the cue and the outcome is sensitive to specific learned cue-reward associations<sup>26,27</sup>. To assess whether the firing of VP neurons is also sensitive to specific cue predictions,



we trained a new cohort of rats to associate one ‘non-specific’ cue with unpredictable sucrose/maltodextrin (similar to the ‘random sucrose/maltodextrin’ task), and two ‘specific’ cues, the first fully predicting sucrose and the second fully predicting maltodextrin, and recorded VP neurons ( $n=487$ ) while they performed the task (Extended Data Figs. 7 and 8). This task is somewhat unusual in that the predicted outcomes are similar in value, perhaps explaining why rats did not noticeably adjust their behavior according to each cue’s prediction (Extended Data Fig. 8b). To quantify how the specific cue predictions modulated outcome-evoked firing rates, the RPE, current outcome and unmodulated models were augmented with two free parameters to estimate the contribution of the new cues; thus, each neuron was fitted with six models. Although we replicated our finding that a subset of VP neurons encodes Rescorla–Wagner RPEs at the time of the outcome (Extended Data Fig. 8g–j), we did not observe many neurons with a TD-like impact of specific reward-predicting cues on their outcome-evoked activity (Extended Data Fig. 8g,k). We did, however, find that the cue-evoked activity of 29% of neurons was impacted by cue identity (Extended Data Fig. 8m). These cells were not more likely to be modulated by cues at the time of the outcome ( $P=0.07$ ). Importantly, these cells tended to have elevated firing for the sucrose cue and reduced firing for the maltodextrin cue relative to the non-specific cue (Extended Data Fig. 8n,o), indicating that the relative values of sucrose and maltodextrin are represented in the firing evoked by their respective predictive cues. Therefore, in this task, where the cues do not overtly influence behavior, VP neurons had cue-evoked, but not outcome-evoked, activity that followed the pattern of a TD error, but additional experiments with more salient cue–reward associations are necessary for definitive conclusions.

## Discussion

We investigated the influence of outcome history on reward-evoked firing in the VP through the lens of RPE signaling. Random presentations of reward revealed a subset of VP neurons that reflected an RPE generated from previously received outcomes and consistent with reward preference. This RPE signal correlated with measures of task engagement in the subsequent trial, and optogenetic manipulation of the VP during reward delivery predictably altered subsequent task engagement. Furthermore, we found that VP RPE neurons demonstrate the expected adaptation when the same reward is presented repeatedly. This series of findings is strong evidence for encoding of outcome history-based RPEs by VP neurons and suggests a role for this signal in adaptive reward seeking.

**An RPE signal beyond dopamine neurons.** A longstanding view is that dopamine neurons compute RPEs locally by integrating distinct elements of the signal relayed from different input regions<sup>4,28,29</sup>. Previous work has revealed that different components of the dopamine neuron RPE calculation depend on various inputs, including lateral habenula<sup>30,31</sup>, rostromedial tegmental nucleus<sup>32,33</sup>, orbitofrontal cortex<sup>22</sup>, ventral striatum<sup>23</sup> and  $\gamma$ -aminobutyric acid (GABA)ergic neurons in the ventral tegmental area (VTA)<sup>27</sup>. A pioneering study on neural activity of monosynaptic inputs to dopamine neurons revealed a mixture of reward and expectation signals across brain regions (including the VP), but notably there were very few upstream neurons encoding full RPEs<sup>14</sup>, maintaining the idea that, by and large, the RPE is calculated within dopamine neurons themselves<sup>29</sup>.

The focus on the construction of an RPE signal within dopamine neurons has left RPE correlates in other reward-processing regions less explored. In the present study, we describe a robust RPE signal in the VP, a region within a highly interconnected circuit implicated in reward learning and reward processing<sup>5,6</sup>. Previously, we characterized a relative value signal in the VP by presenting rats with various combinations of differentially preferred rewards<sup>9</sup>. In that work, we observed an influence of only one previous trial on

reward-evoked signaling when looking at the full recorded population, a result we replicated here by performing an outcome history regression on all VP neurons (Fig. 1g). Our innovation in the present work is implementing a computational modeling approach that allowed us to identify individual neurons with firing consistent with RPEs, integrating outcome history over several trials.

There have been a few previous attempts to characterize RPE-like signals in the VP<sup>12,14–16</sup>, with mixed results. An important distinction in our present work is that we focused most of our analysis on Rescorla–Wagner trial-based RPEs, which integrate over outcome history, whereas previous studies looked for outcome signaling modulation by specific predictive cues within a TD learning framework. Both models have been useful for characterizing dopamine neuron activity. Our data here indicate that, much like dopamine neurons<sup>4,20,25</sup>, a subset of VP neurons encode trial-based RPEs. In addition, we linked VP activity during the reward epoch with changes in task engagement, indicating that this signal may contribute to updating the estimate of the task’s value, consistent with the Rescorla–Wagner model. We once again found mixed evidence for TD error signals in VP.

One possible explanation for the apparent lack of TD signaling is that the VP may not update the values of particular cues, but rather may update the estimate of average environmental reward over behaviorally relevant timescales. Theories and experiments have suggested that average environmental reward signals are critical for invigorating behavior<sup>34–36</sup>. Intriguingly, subtle manipulations of VP slow response vigor<sup>8</sup> and gross manipulations are typically associated with motivational deficits<sup>5,6</sup>. Both of these effects are consistent with a role for the VP in computing average reward. Our finding that VP activity correlates with subsequent task engagement and that VP optogenetic manipulations alter subsequent task engagement additionally supports this idea. As the cue–response contingency was identical on all trial types in the tasks presented in this study, future work will need to clarify whether this signal updates estimates of global reward rate, or perhaps the value of specific reward-seeking actions. In addition, a task with greater motivation and learning demands could help distinguish the roles of RPE and current outcome cells.

**Relationship between the VP and dopamine RPE signals.** VP neurons have direct and indirect reciprocal connections with VTA dopamine neurons, so a natural question is how RPE signals in each population may influence each other. Building on the discussion in An RPE signal beyond dopamine neurons, it is possible that VP integrates reward history to provide an average-reward error signal to dopamine neurons. Indeed, dopamine activity tracks average reward<sup>25,34,35</sup>. As we saw less RPE-like activity in the NAc, another input to midbrain dopamine neurons, the VP could be a privileged source for this information. There are multiple routes for VP activity to reach the VTA, given the demonstrations of VP synapses not only on to VTA neurons<sup>6,14,18,37</sup> but also on to VTA input nuclei such as the lateral habenula and rostromedial tegmental nucleus<sup>17,18,38,39</sup>. Stimulation of VP GABAergic neurons increases the number of putative midbrain dopamine neurons expressing Fos, consistent with an indirect mechanism for modulating features of dopamine neuron RPE signaling<sup>18</sup>. It is interesting that, in songbirds, the VP has been shown to send performance-related error signals to the VTA during singing<sup>40–42</sup>.

On the other hand, the VTA could be a source for VP RPE signals; in addition to dopaminergic innervation of the VP<sup>5,6,37</sup>, the VTA sends dense glutamatergic projections, which may be more likely to mediate the phasic responses we observed in VP<sup>43</sup>. Future work should untangle a uni- or bidirectional influence of error-related signals in these regions, as well as possible unique roles of each population in adaptive behavior. A combination of projection-specific recordings and manipulations in the VP and

VTA would help clarify this question. In addition, other candidate downstream regions whereby VP RPE signals may mediate learning and behavior include the mediodorsal nucleus (the classic VP thalamic output)<sup>5–7,44</sup>, as well as the lateral hypothalamus<sup>45</sup> and lateral habenula<sup>17,18,38,39</sup>.

One metric on which signaling in the VP and VTA can be compared is the prevalence of outcome history-sensitive RPEs across the population. In our task, we found as many as ~30% of neurons encoded RPEs, but this was variable across tasks. This proportion is similar to that found in dopamine neurons, which ranges from 15% to 50% in rodents<sup>20–22</sup>. In our dataset, we noted that the task with the greatest range in reward value revealed the most RPE cells. Future work should explore how additional changes in task parameters, such as the inclusion of aversive outcomes<sup>14,15</sup>, Pavlovian versus instrumental contingencies<sup>46,47</sup> and deterministic versus probabilistic outcomes, impact RPE signaling in the VP relative to the NAC.

### Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41593-020-0688-5>.

Received: 16 October 2019; Accepted: 7 July 2020;

Published online: 10 August 2020

### References

- Sutton, R. S. & Barto, A. G. *Introduction to Reinforcement Learning* (MIT Press, Cambridge, MA, 1998).
- Rescorla, R. A. & Wagner, A. R. A theory of Pavlovian conditioning: variations in the effectiveness of reinforcement and nonreinforcement, in *Classical Conditioning II: Current Research and Theory*, Vol. 2 (eds Black, A. H. & Prokasy, W. F.), 64–99 (Apple-Century-Crofts, 1972).
- Schultz, W., Dayan, P. & Montague, P. R. A neural substrate of prediction and reward. *Science* **275**, 1593–1599 (1997).
- Bayer, H. M. & Glimcher, P. W. Midbrain dopamine neurons encode a quantitative reward prediction error signal. *Neuron* **47**, 129–141 (2005).
- Smith, K. S., Tindell, A. J., Aldridge, J. W. & Berridge, K. C. Ventral pallidum roles in reward and motivation. *Behav. Brain Res.* **196**, 155–167 (2009).
- Root, D. H., Melendez, R. I., Zaborszky, L. & Napier, T. C. The ventral pallidum: subregion-specific functional anatomy and roles in motivated behaviors. *Prog. Neurobiol.* **130**, 29–70 (2015).
- de Olmos, J. S. & Heimer, L. The concepts of the ventral striatopallidal system and extended amygdala. *Ann. NY Acad. Sci.* **877**, 1–32 (1999).
- Richard, J. M., Ambroggi, E., Janak, P. H. & Fields, H. L. Ventral pallidum neurons encode incentive value and promote cue-elicited instrumental actions. *Neuron* **90**, 1165–1173 (2016).
- Ottenheimer, D., Richard, J. M. & Janak, P. H. Ventral pallidum encodes relative reward value earlier and more robustly than nucleus accumbens. *Nat. Commun.* **9**, 4350 (2018).
- Fujimoto, A. et al. Signaling incentive and drive in the primate ventral pallidum for motivational control of goal-directed action. *J. Neurosci.* **39**, 1793–1804 (2019).
- White, J. K. et al. A neural network for information seeking. *Nat. Commun.* **10**, 1–19 (2019).
- Tindell, A. J., Berridge, K. C. & Aldridge, J. W. Ventral pallidal representation of Pavlovian cues and reward: population and rate codes. *J. Neurosci.* **24**, 1058–1069 (2004).
- Tachibana, Y. & Hikosaka, O. The primate ventral pallidum encodes expected reward value and regulates motor action. *Neuron* **76**, 826–837 (2012).
- Tian, J. et al. Distributed and mixed information in monosynaptic inputs to dopamine neurons. *Neuron* **91**, 1374–1389 (2016).
- Stephenson-Jones, M. et al. Opposing contributions of gabaergic and glutamatergic ventral pallidal neurons to motivational behaviors. *Neuron* **105**, 921–933 (2020).
- Kaplan, A., Mizrahi-Kliger, A. D., Israel, Z., Adler, A. & Bergman, H. Dissociable roles of ventral pallidum neurons in the basal ganglia reinforcement learning network. *Nat. Neurosci.* **23**, 556–564 (2020).
- Tooley, J. et al. Glutamatergic ventral pallidal neurons modulate activity of the habenula–tegmental circuitry and constrain reward seeking. *Biol. Psychiatry* **83**, 1012–1023 (2018).
- Faget, L. et al. Opponent control of behavioral reinforcement by inhibitory and excitatory projections from the ventral pallidum. *Nat. Commun.* **9**, 849 (2018).
- Sclafani, A., Hertwig, H., Vigorito, M. & Feigin, M. B. Sex differences in polysaccharide and sugar preferences in rats. *Neurosci. Biobehav. Rev.* **11**, 241–251 (1987).
- Mohebi, A. et al. Dissociable dopamine dynamics for learning and motivation. *Nature* **570**, 65–70 (2019).
- Roesch, M. R., Calu, D. J. & Schoenbaum, G. Dopamine neurons encode the better option in rats deciding between differently delayed or sized rewards. *Nat. Neurosci.* **10**, 1615 (2007).
- Takahashi, Y. K. et al. Expectancy-related changes in firing of dopamine neurons depend on orbitofrontal cortex. *Nat. Neurosci.* **14**, 1590 (2011).
- Takahashi, Y. K., Langdon, A. J., Niv, Y. & Schoenbaum, G. Temporal specificity of reward prediction errors signaled by putative dopamine neurons in rat VTA depends on ventral striatum. *Neuron* **91**, 182–193 (2016).
- Sutton, R. S. Learning to predict by the methods of temporal differences. *Mach. Learn.* **3**, 9–44 (1988).
- Nakahara, H., Itoh, H., Kawagoe, R., Takikawa, Y. & Hikosaka, O. Dopamine neurons can represent context-dependent prediction error. *Neuron* **41**, 269–280 (2004).
- Fiorillo, C. D., Tobler, P. N. & Schultz, W. Discrete coding of reward probability and uncertainty by dopamine neurons. *Science* **299**, 1898–1902 (2003).
- Eshel, N. et al. Arithmetic and local circuitry underlying dopamine prediction errors. *Nature* **525**, 243–246 (2015).
- Keiflin, R. & Janak, P. H. Dopamine prediction errors in reward learning and addiction: from theory to neural circuitry. *Neuron* **88**, 247–263 (2015).
- Watabe-Uchida, M., Eshel, N. & Uchida, N. Neural circuitry of reward prediction error. *Annu. Rev. Neurosci.* **40**, 373–394 (2017).
- Matsumoto, M. & Hikosaka, O. Lateral habenula as a source of negative reward signals in dopamine neurons. *Nature* **447**, 1111–1115 (2007).
- Tian, J. & Uchida, N. Habenula lesions reveal that multiple mechanisms underlie dopamine prediction errors. *Neuron* **87**, 1304–1316 (2015).
- Jhou, T. C., Fields, H. L., Baxter, M. G., Saper, C. B. & Holland, P. C. The rostromedial tegmental nucleus (RMTg), a GABAergic afferent to midbrain dopamine neurons, encodes aversive stimuli and inhibits motor responses. *Neuron* **61**, 786–800 (2009).
- Hong, S., Jhou, T. C., Smith, M., Saleem, K. S. & Hikosaka, O. Negative reward signals from the lateral habenula to dopamine neurons are mediated by rostromedial tegmental nucleus in primates. *J. Neurosci.* **31**, 11457–11471 (2011).
- Niv, Y., Daw, N. D., Joel, D. & Dayan, P. Tonic dopamine: opportunity costs and the control of response vigor. *Psychopharmacology* **191**, 507–520 (2007).
- Hamid, A. A. et al. Mesolimbic dopamine signals the value of work. *Nat. Neurosci.* **19**, 117–126 (2016).
- Bari, B. A. et al. Stable representations of decision variables for flexible behavior. *Neuron* **103**, 922–933 (2019).
- Beier, K. T. et al. Circuit architecture of vta dopamine neurons revealed by systematic input–output mapping. *Cell* **162**, 622–634 (2015).
- Hong, S. & Hikosaka, O. Diverse sources of reward value signals in the basal ganglia nuclei transmitted to the lateral habenula in the monkey. *Front. Hum. Neurosci.* **7**, 778 (2013).
- Knowland, D. et al. Distinct ventral pallidal neural populations mediate separate symptoms of depression. *Cell* **170**, 284–297 (2017).
- Gale, S. D. & Perkel, D. J. A basal ganglia pathway drives selective auditory responses in songbird dopaminergic neurons via disinhibition. *J. Neurosci.* **30**, 1027–1037 (2010).
- Chen, R. et al. Songbird ventral pallidum sends diverse performance error signals to dopaminergic midbrain. *Neuron* **103**, 266–276 (2019).
- Kearney, M. G., Warren, T. L., Hissey, E., Qi, J. & Mooney, R. Discrete evaluative and premotor circuits enable vocal learning in songbirds. *Neuron* **104**, 559–575 (2019).
- Hnasko, T. S., Hjelmstad, G. O., Fields, H. L. & Edwards, R. H. Ventral tegmental area glutamate neurons: electrophysiological properties and projections. *J. Neurosci.* **32**, 15076–15085 (2012).
- Leung, B. K. & Balleine, B. W. Ventral pallidal projections to mediodorsal thalamus and ventral tegmental area play distinct roles in outcome-specific Pavlovian-instrumental transfer. *J. Neurosci.* **35**, 4953–4964 (2015).
- Prasad, A. A. et al. Complementary roles for ventral pallidum cell types and their projections in relapse. *J. Neurosci.* **40**, 880–893 (2020).
- Richard, J. M., Stout, N., Acs, D. & Janak, P. H. Ventral pallidal encoding of reward-seeking behavior depends on the underlying associative structure. *eLife* **7**, e33107 (2018).
- Ottenheimer, D. J., Wang, K., Haimbaugh, A., Janak, P. H. & Richard, J. M. Recruitment and disruption of ventral pallidal cue encoding during alcohol seeking. *Eur. J. Neurosci.* **50**, 3428–3444 (2019).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

© The Author(s), under exclusive licence to Springer Nature America, Inc. 2020

## Methods

**Animals.** Subjects for electrophysiology experiments were male Long–Evans rats ( $n = 15$ ) from Envigo weighing 250–275 g and aged 60 d on arrival. Subjects for the optogenetic experiment were male ( $n = 14$ ) and female ( $n = 17$ ) Long–Evans rats from Envigo weighing 200–275 g and aged 60 d on arrival. Rats were single housed on a 12-h light:dark cycle and given free access to food and water in their home cages for the duration of the experiment. All experimental procedures were performed in strict accordance with protocols approved by the Animal Care and Use Committee at Johns Hopkins University.

**Reward solutions.** We used 10% solutions by weight of sucrose (Thermo Fisher Scientific) and maltodextrin (SolCarb, Solace Nutrition) in tap water (or tap water alone). Rats were given free access to the solutions in their home cages before training began to ensure familiarity with the rewards.

**Behavioral tasks.** *Random sucrose/maltodextrin.* Data from this task have been previously published<sup>3</sup>. Rats ( $n = 11$ ) were trained to respond to a 10-s white noise cue by making an entry into the reward port. The cue terminated on port entry and, 500 ms after port entry, 110  $\mu$ l of either reward was delivered into the metal cup within the reward port over 2 s. Sucrose and maltodextrin trials were pseudorandomly interspersed throughout the session, such that rats could not detect the identity of the reward until it was delivered. Individual licks were recorded with a custom-built, arduino-based lickometer using a capacitance sensor (MPRI21, Adafrit Industries) with a 1-kHz sampling rate. Each cue was separated by a variable ITI that averaged 45 s. During the ITI, the reward cup was evacuated via a vacuum pump, flushed with 110  $\mu$ l of water and evacuated again. Maltodextrin, sucrose and water were each delivered via separate infusion pumps (Med Associates) and separate metal tubes entering the cup. There were 60 trials per session. For 3 rats, 2 additional sessions were conducted with tap water as a third outcome for a total of 90 trials (random sucrose/maltodextrin/water).

*Blocked sucrose/maltodextrin.* For the same group of rats as the random task, additional blocked sessions were performed on alternating days with the random sucrose/maltodextrin task. In blocked sessions, sucrose and maltodextrin were presented 30 trials in a row for a total of 60 trials. The order of the rewards switched each blocked session. Predictable and random sucrose/maltodextrin: a new group of rats ( $n = 4$ ) was trained on a task with the same trial structure (with a shorter ITI of 30 s) but with three possible auditory cues. One predicted sucrose delivery with 100% probability (30 trials), one maltodextrin with 100% probability (30 trials) and one, as in the random sucrose/maltodextrin task, predicted each reward with a 50% probability (60 trials).

**Preference test.** To assay rats' preference for sucrose or maltodextrin, we performed two 60-min, two-bottle choice tests, during which rats had free access to 10% solutions of each reward. Bottles were weighed before and after to determine the amount of each solution consumed by each rat. The first test was after recovery from surgery and before recording. The second was at least a day after the final session with sucrose and maltodextrin.

**Surgical procedures.** Rats were anesthetized with isoflurane (5%) and maintained under anesthesia for the duration of the surgery (1–2%). Rats received injections of carprofen (5 mg kg<sup>-1</sup>) and cefazolin (70 mg kg<sup>-1</sup>) before incision.

**Electrophysiology studies.** Drivable electrode arrays were prepared with three-dimensional-printed plastic housing, 16 insulated tungsten wires and 2 silver ground wires. The drives were surgically implanted in trained rats. Using a stereotaxic arm, electrodes were aimed at the VP ( $n = 9$ , anteroposterior (AP) +0.5 mm, mediolateral (ML) +2.4 mm, dorsoventral (DV) –8 mm) or NAc ( $n = 6$ , AP +1.5 mm, ML +1.2 mm, DV –7 mm).

**Optogenetic studies.** First, 0.7  $\mu$ l of virus containing the archaerhodopsin gene construct (AAV5-CamKIIa-eArchT3.0-eYFP,  $7 \times 10^{12}$  viral particles per ml from the University of North Carolina Vector Core), ChR2 (AAV5-hsyn-hChR2(H134R)-eYFP,  $1.7 \times 10^{13}$  viral particles per ml from Addgene, gift from K. Deisseroth) or their respective control virus (AAV5-CamKIIa-eYFP,  $7.4 \times 10^{12}$  viral particles per ml from the University of North Carolina Vector Core, or AAV5-hsyn-EGFP,  $1.2 \times 10^{13}$  viral particles per ml from Addgene, gift from B. Roth) was delivered bilaterally to the VP through 31-G, gas-tight Hamilton syringes at a rate of 0.1  $\mu$ l min<sup>-1</sup> for 7 min controlled by a Micro4 Ultra Microsyringe Pump 3 (World Precision Instruments). Injectors were left in place for 10 min after the infusion to allow virus to diffuse away from the infusion site. Injector tips were aimed at the following coordinates in relation to bregma: +0.5 mm AP,  $\pm 2.5$  mm ML, –8.2 mm DV. Then, rats were implanted with 300- $\mu$ m-diameter optic fibers constructed in house, aimed 0.3 mm above the center of the virus infusion.

**Electrophysiological recording.** After recovery in their home cages, rats were trained on the task again, until they became accustomed to the recording setup. Electrical signals and behavioral events were collected using the OmniPlex

system (Plexon) with a 40-kHz sampling rate. For rats in the random and blocked sucrose/maltodextrin tasks, we continued to record from the same location for multiple sessions if new neurons appeared on previously unrecorded channels. For the random task, if multiple sessions from the same location were included in the analysis, the same wire was never included more than once. For the blocked task, we occasionally included a wire in the same location twice if each of the two sessions had a different block order. If no neurons were detectable or after successful recording, the drive was advanced 160  $\mu$ m, and recording resumed in the new location at a minimum 2 d later to ensure settling of the tissue around the wires. For rats in the predictable and random sucrose/maltodextrin tasks, we maintained the wires in the same position for the duration of the experiment. Each wire from these rats contributed to the included dataset only once.

**Optogenetic manipulations.** *Inhibition.* At least 5 weeks after surgery and completion of operant training, rats were habituated to patch-cord connections. Animals were connected via a ceramic mating sleeve to a 200- $\mu$ m-core patch cord, which was then connected through a fiberoptic rotary joint (Doric), to another patch cord that interfaced with a 532-nm DPSS laser (Opto-Engine LLC). The time of laser delivery was initiated by transistor–transistor logic pulses from MedPC SmartCTRL cards to a Master9 Stimulus Controller (AMPI), which dictated the duration of stimulation. For this experiment, rats were trained on a variation of the random sucrose/maltodextrin task in which, rather than maltodextrin delivery, rats received sucrose + continuous (5 s, 15–20 mW) photoinhibition of the VP. We chose 5 s to span the majority of consumption (Extended Data Fig. 4a) and minimize the possibility that the release of inhibition within the first few seconds would also impact the behavior, because reward-specific signaling persists for up to 4 s (Fig. 1e). For these sessions the reward volume was reduced to 55  $\mu$ l and the total number of trials was increased to 90. In our analysis, we included only rats that completed at least 30 trials and had both fibers and viral expression in the VP. This resulted in seven rats in each group: four males and three females in the ArchT3.0 group, and two males and five females in the YFP group.

*Excitation.* we used the same protocol as the inhibition group, but with unilateral 40-Hz pulsed photoexcitation of the VP for 2 s (10-ms pulse width, 10–12 mW). We also conducted a session in which stimulation occurred at cue onset for 2 s (or until reward delivery if sooner). As rats were implanted bilaterally, we stimulated the side with maximal effect and minimal off-target effects as determined in previous experiments. We included only rats that completed at least 30 trials and had their stimulated fiber and viral expression in the VP. This resulted in five males and five females in the ChR2 group, and three males and four females in the GFP group.

**Histology.** Animals were deeply anesthetized with pentobarbital. For rats in the electrophysiology experiments, electrode sites were labeled by passing a DC current through each electrode. All rats were perfused intracardially with 0.9% saline followed by 4% paraformaldehyde, post-fixed and sectioned into 50- $\mu$ m slices on a cryostat. Cresyl violet (electrophysiology experiments) or DAPI (optogenetic experiments) was used to visualize electrode, virus and fiber placements.

**Spike sorting and initial analysis.** Spikes were sorted into units using offline sorter (Plexon); after initial manual selection of units, based on clustering of waveforms along the first two principal components, units were separated and refined using waveform energy and waveform heights at various times relative to threshold crossing (slices). Any units that were not detectable for the entire session were discarded. Event creation and review of individual neurons' responses were conducted in NeuroExplorer (Nex Technologies). Cross-correlation was plotted for simultaneously recorded units to identify and remove any neurons that were recorded on multiple channels. All subsequent analysis was performed in MATLAB (MathWorks).

**PSTH creation.** PSTHs were constructed using 0.01-ms bins surrounding the event of interest (generally, reward delivery). They were smoothed using a half-normal filter ( $\sigma = 6.6$ ) that used only activity in previous, but not upcoming, bins. Each bin of the PSTH was z-scored by subtracting the mean firing rate across 10-s windows before each trial and dividing by the s.d. across those windows ( $n = \text{number of trials}$ ). PSTHs for licking were created in the same manner (without z-scoring) using 0.05-ms bins and  $\sigma = 8$ . For the neural activity plots with trials binned by RPE, we smoothed both individual trials ( $\sigma = 10$ ) and the entire PSTH ( $\sigma = 20$ ) with half-normal filters, because each trace was constructed from just a handful of trials.

**Model fitting.** For each neuron, we took the spike count,  $s(t)$ , within the 0.75- to 1.95-s post-reward delivery time bin for each trial and fit Poisson's spike count models using maximum likelihood estimation, which is a method for estimating the parameters that best predict the observed data (in this case, spike counts), under the assumed model. For the random and blocked sucrose/maltodextrin tasks, we fit the following three models.



RPE model:

$$\begin{aligned}\delta(t) &= o(t) - V(t) \\ V(t+1) &= V(t) + \alpha \cdot \delta(t) \\ s(t) &\sim \text{Poisson}(\exp(\alpha \cdot \delta(t) + b))\end{aligned}$$

where  $V(t)$  is the expected value,  $\delta(t)$  is the RPE,  $o(t)$  is the outcome and  $\alpha$  is the learning rate. For the tasks with sucrose and maltodextrin outcomes, we coded  $o(t)=0$  for maltodextrin and 1 for sucrose. For the tasks with sucrose, maltodextrin and water outcomes, we coded  $o(t)=0$  for water, 1 for sucrose and  $\rho$  for maltodextrin, a free parameter we estimated during model fitting. To map RPEs to spike counts, we used  $a$  as a slope (gain) and  $b$  as an intercept (offset) parameter. This affine-transformed RPE was mapped through an exponential function, to avoid negative values, and used as the rate parameter for a Poisson distribution. To identify neurons with value responses at the time of the cue, we replaced  $\delta(t)$  with  $V(t)$  in Poisson's function mapping latent variables to spike counts.

Current outcome model:

$$s(t) \sim \text{Poisson}(\exp(a \cdot o(t) + b))$$

Unmodulated model:

$$s(t) \sim \text{Poisson}(\exp(\ln(\bar{s})))$$

where  $\bar{s}$  is the mean firing rate.

For the predictable and random sucrose/maltodextrin task, we added the following three models:

RPE + cue model:

$$\begin{aligned}\delta(t) &= o(t) - V(t) \\ V(t+1) &= V(t) + \alpha \cdot \delta(t)\end{aligned}$$

If a sucrose-predicting cue was given:

$$s(t) \sim \text{Poisson}(\exp(a \cdot \delta(t) + b - V_{\text{sucrose}}))$$

If a maltodextrin-predicting cue was given:

$$s(t) \sim \text{Poisson}(\exp(a \cdot \delta(t) + b - V_{\text{maltodextrin}}))$$

If a nonpredictive cue was given:

$$s(t) \sim \text{Poisson}(\exp(a \cdot \delta(t) + b))$$

where  $V_{\text{sucrose}}$  and  $V_{\text{maltodextrin}}$  are free parameters for the values of the sucrose- and maltodextrin-predicting cues, respectively.

Current outcome + cue model:

If a sucrose-predicting cue was given:

$$s(t) \sim \text{Poisson}(\exp(a \cdot o(t) + b - V_{\text{sucrose}}))$$

If a maltodextrin-predicting cue was given:

$$s(t) \sim \text{Poisson}(\exp(a \cdot o(t) + b - V_{\text{maltodextrin}}))$$

If a nonpredictive cue was given:

$$s(t) \sim \text{Poisson}(\exp(a \cdot o(t) + b))$$

Unmodulated + cue model:

If a sucrose-predicting cue was given:

$$s(t) \sim \text{Poisson}(\exp(\ln(\bar{s}) - V_{\text{sucrose}}))$$

If a maltodextrin-predicting cue was given:

$$s(t) \sim \text{Poisson}(\exp(\ln(\bar{s}) - V_{\text{maltodextrin}}))$$

If a nonpredictive cue was given:

$$s(t) \sim \text{Poisson}(\exp(\ln(\bar{s})))$$

To estimate predictive cue effects on firing at the time of the cue, we fit the unmodulated model and unmodulated + cue model, with  $V_{\text{sucrose}}$  and  $V_{\text{maltodextrin}}$  sign-flipped.

We also considered RPE models in which the predictive cue allowed for partial to full cancellation of RPEs.

If a sucrose-predicting cue was given:

$$\eta(t) = o(t) - ((1-w) \cdot V(t) + w \cdot V_{\text{sucrose}})$$

If a maltodextrin-predicting cue was given:

$$\eta(t) = o(t) - ((1-w) \cdot V(t) + w \cdot V_{\text{maltodextrin}})$$

If a nonpredictive cue was given:

$$\eta(t) = o(t) - V(t)$$

$$s(t) \sim \text{Poisson}(\exp(a \cdot \eta(t) + b))$$

We fixed  $V_{\text{sucrose}} = 1$  and  $V_{\text{maltodextrin}} = 0$  and set  $w$  as a free parameter. If  $w=0$ , this is equivalent to the RPE model and, if  $w=1$ , the predictive cues allow for full cancellation of the RPE ( $\eta(t)=0$ ). Intermediate values of  $w$  allow the predictive cues to partially cancel the outcome history-based RPE. This model was best for a negligible number of neurons.

We analyzed only trials in which the rat licked within the first 2 s of reward delivery, to ensure that it sampled the outcome. For all RPE models except those in the three-outcome task,  $V(1)$  was initialized to 0.5. For the three-outcome task,  $V(1)$  was initialized to  $\frac{1+\rho}{3}$ , where  $\rho$  is the estimated value of maltodextrin, to avoid biasing the initial value estimation. For all models with a slope parameter, we constrained the slope,  $a$ , to be  $>0$ , as previous work showed that a trivial fraction of VP neurons preferentially encodes low-value rewards (9). We found maximum likelihood estimates for each model and selected the best model using an AIC (a lower AIC indicates a better fit, after taking into account the number of parameters). We used 10 randomly selected starting initial values for each parameter to avoid finding local minima.

We performed a number of checks to ensure reliability of the model classification approach. First, we conducted a model recovery exercise by generating simulated neurons for each of our three models (see "Model recovery"). We then classified these simulated neurons (blinded to the true category) and asked how many were correctly recovered. Overall, simulated neurons were classified as the correct model most of the time (Extended Data Fig. 2c). Moreover, our classification of RPE neurons was probably on the conservative end; RPE neurons were more often classified as current outcome or unmodulated neurons than vice versa. We also found that our approach allowed for unbiased estimates of the parameters used to simulate the neurons (Extended Data Fig. 2d). Second, as an additional test of reliability, we repeated the classification of VP and NAc neurons with the Bayesian information criterion (BIC), which more harshly punishes additional free parameters. Although this approach classified fewer neurons as RPE and current outcome in both the VP and the NAc, the pattern of results remained the same, with more RPE and current outcome cells in the VP than in the NAc, and a higher correlation between RPE model predictions and actual spike counts in the VP than in the NAc (Extended Data Fig. 9a–e). We further analyzed the subset of VP neurons classified as RPE cells by AIC but not by BIC. Intuitively, as this subset of neurons less robustly encodes RPEs, we sought to determine whether they were simply noise or probably encoding RPEs. These AIC RPE cells were strongly modulated by the previous outcome and demonstrated firing patterns typical of the broader RPE population (Extended Data Fig. 9g–i), suggesting that RPE was the appropriate classification and BIC is probably too conservative for estimating the prevalence of outcome history-sensitive signaling across the population. Therefore, we used AIC for all analyses in the main text and figures.

We also considered a number of alternatives to our standard RPE model. First, we fit asymmetric learning models, with separate learning rates for positive prediction errors and negative prediction errors ( $\alpha_{\text{PPE}}$  and  $\alpha_{\text{NPE}}$ , respectively). These models allow for biased (optimistic or pessimistic) estimates of the value function. The asymmetric learning models were initialized with  $V(1) = \frac{\alpha_{\text{PPE}}}{\alpha_{\text{PPE}} + \alpha_{\text{NPE}}}$ , the steady-state-biased value estimate. Note that if  $\alpha_{\text{PPE}} = \alpha_{\text{NPE}}$ , then  $V(1) = 0.5$ , the initialization we used for our single learning-rate models. These models fit best for a small number of neurons (of 75 total RPE neurons in the 'random sucrose/maltodextrin' task, 15 were best fit by the asymmetric RPE model), suggesting that asymmetric RPE coding is not a cardinal feature.

Second, we allowed the slope parameter,  $a$ , to be negative, to capture neurons that signal negative RPEs or negative outcomes with an increase in firing rate. Only 5 of 77 RPE neurons and 16 of 142 current outcome neurons were best fit by negative-slope models, suggesting that the vast majority of neurons in the VP signal positive RPEs or positive outcomes with an increase in firing rate.

Third, we considered two phenomenological models to explain the firing rates of VP neurons. In both models, spike counts were a function of the reward (that is,  $s(t) \approx \text{Poisson}(\exp(a \times o(t) + b))$ ). The first, the 'habituation' model, allowed the slope parameter,  $a$ , to vary as a function of recent reward history. We allowed the slope to decrease closer to 0 when reward history was greater (sucrose in the recent past), and increase when reward history was poorer. We recoded  $o(t)$  to 0.5 for sucrose and  $-0.5$  for maltodextrin, to allow the slope to modulate the firing rate on maltodextrin trials (if  $o(t)=0$  on maltodextrin trials, then  $a \times o(t)=0$ , regardless of the slope). This allowed spike counts to decrease when sucrose or maltodextrin was presented repeatedly, allowing for habituation. The second, the 'adaptation' model, allowed the intercept parameter,  $b$ , to vary as a function of recent reward history. This allowed spike counts to decrease suddenly when maltodextrin was delivered after a string of sucrose rewards, and increase suddenly when sucrose was delivered after a string of maltodextrin rewards, simulating RPE-like signals. When we fit all models to firing rates (RPE, current outcome, unmodulated, habituation, adaptation), the habituation model fit best for 33 neurons and the adaptation model for 18 neurons. The RPE model fit best for 54 neurons, which was significantly more than the habituation model ( $P=0.02$ ) and the adaptation model ( $P=9 \times 10^{-6}$ ).

Fourth, to test whether our procedure for classifying neurons was invariant to transformations of the data, we z-scored the spike counts and fit Gaussian observation models, with an extra parameter to estimate the variance. Using these



models, we identified 66 of 436 (15%) neurons as RPE coding, similar to the 72 of 436 (17%) we identified with Poisson's observation model.

Fifth, a recent study has argued that neurons may be inappropriately classified as coding a latent variable if there are temporal correlations in both the latent variable and the firing rates<sup>48</sup>. This concern does not hold for our study because, in all but the blocked task, the reward was delivered randomly. This ensures that the estimated RPE signal is not autocorrelated. However, to test this concern rigorously, we adapted and implemented the permutation test recommended by the authors of the article<sup>48</sup>. The permutation test identifies neurons that are more correlated with RPEs estimated from that session than from other sessions. As such, it is an exceedingly strict test. We first fit a linear regression to each neuron to predict z-scored spike count as a function of RPEs. The resulting *t*-statistic was then compared with a null distribution of *t*-statistics, generated by fitting similar linear regressions using RPEs from all other sessions. A neuron was considered to encode RPEs if the *t*-statistic fell outside the 5% significance boundary of the null distribution. We found that 91% of our model-identified RPE neurons were considered to be RPE neurons by the permutation test. This was probably because the correlation between RPEs across the session was not different from chance (median correlation between RPEs ( $\pm 95\%$  bootstrapped CI: 0.112 (0.107–0.117); median correlation between two random Gaussian vectors: 0.105 (0.103–0.108)).

**Correlation and RPE tuning curves for real and simulated neurons.** For neurons best fit by the RPE model, we report correlations between real and predicted spike trains, as well as RPE tuning curves for real and predicted spike counts. For each neuron, we estimated Pearson's correlation coefficient between real spike counts and 501 independent model-generated spike count trains, using parameters estimated from the same neuron, and report the median correlation. The median-correlated spike count is plotted in Figs. 2e and 3i. To generate RPE tuning curves for real spike counts, we took z-scored spike counts and binned according to the estimated RPEs. We performed this procedure for all RPE neurons and report the average tuning curve. To generate tuning curves for predicted spike counts, we simulated spike trains using neuron-derived parameter estimates and followed the same procedure.

**Model recovery.** We simulated 200 RPE model neurons, 200 current outcome model neurons and 200 unmodulated model neurons to assess whether our modeling recovery strategy could correctly classify neurons. For each neuron, we simulated 55 trials of the random sucrose/maltodextrin task. We constrained  $\alpha$  to 0 to 1, the slope (*a*) to 1 to 4 and the intercept (*b*) to  $-5$  to  $5$ . We again used 10 randomly selected starting initial values for each parameter to avoid finding local minima.

**Outcome history-based linear regression.** To estimate how the outcome of the current and previous trials affected the firing rate of the current trial, we conducted a complete-pooling linear regression analysis. We z-scored the firing rate of each neuron using the baseline activity across the set of 10-s bins before each trial and combined the firing rates of all neurons of interest. Similarly, our design matrix included the current and 10 previous trial outcomes for all neurons of interest. Significance for each trial lag was determined during model fitting using the *t*-statistic against the null hypothesis that the coefficient for that trial was zero ( $P = 0.05$  cutoff). For the random sucrose/maltodextrin task, we gave maltodextrin a value of 0 and sucrose a value of 1. For the random sucrose/maltodextrin/water task, water was given a value of 0, sucrose a value of 1, and maltodextrin a value of 0.75 for RPE cells and 0.8 for current outcome cells, the values that achieved the maximum  $R^2$  for the linear regression. We followed the same process to generate outcome history regression coefficients for simulated neurons (Extended Data Fig. 2).

As the linear regression model can capture arbitrary linear relationships between outcome history and firing rate, it is not immediately clear what pattern of coefficients should be taken as evidence of RPE coding. We can directly relate the linear regression model to the Rescorla–Wagner model to gain insight. To begin, the regression model takes the following form:

$$f(t) = \beta_{\text{int}} + \beta_0 o(t) + \beta_1 o(t-1) + \beta_2 o(t-2) + \dots + \beta_N o(t-N)$$

where  $f(t)$  is the firing rate on trial *t*,  $o(t)$  is the outcome (1 for sucrose, 0 for maltodextrin),  $\beta_{\text{int}}$  is the regression coefficient for the intercept and  $\beta_i$  is the regression coefficient for  $o(t-i)$ . This can be rearranged as:

$$\frac{f(t) - \beta_{\text{int}}}{\beta_0} = o(t) + \sum_{i=1}^N \frac{\beta_i}{\beta_0} o(t-i)$$

Recall that  $\delta(t) = o(t) - V(t)$  in the Rescorla–Wagner model. We can relate the above equation to  $\delta(t)$  if we assume  $V(t) = -\sum_{i=1}^N \frac{\beta_i}{\beta_0} o(t-i)$ .

$$\frac{f(t) - \beta_{\text{int}}}{\beta_0} = o(t) - V(t) = \delta(t)$$

Therefore, the firing rate (after subtracting  $\beta_{\text{int}}$  and scaling by  $\beta_0$ ) can be linearly related to the RPE. To understand under what conditions we can relate the firing

rate to RPE coding, we can recursively expand  $V(t)$  as follows, where  $\alpha$  is the learning rate<sup>1</sup>:

$$V(t) = (1 - \alpha)^t V(1) + \sum_{i=1}^{t-1} \alpha (1 - \alpha)^{t-i} o(i)$$

allowing us to relate  $V(t)$  in the following two ways (the first from recursively expanding  $V(t)$  above, the second from our assumption that

$$V(t) = -\sum_{i=1}^N \frac{\beta_i}{\beta_0} o(t-i)$$

$$V(t) = \alpha o(t-1) + \alpha(1-\alpha)o(t-2) + \alpha(1-\alpha)^2 o(t-3) + \dots$$

$$V(t) = \frac{-\beta_1}{\beta_0} o(t-1) + \frac{-\beta_2}{\beta_0} o(t-2) + \frac{-\beta_3}{\beta_0} o(t-3) + \dots$$

This means that the regressors ( $\beta_1, \beta_2, \dots$ ) should decay exponentially ( $\frac{-\beta_1}{\beta_0} = \alpha, \frac{-\beta_2}{\beta_0} = \alpha(1-\alpha)$  and so on), and that  $\frac{-\beta_i}{\beta_0} \geq 0$  for  $i \geq 1$ . In summary, regression coefficients with a positive  $\beta_0$  and negative  $\beta_i$  values (for  $i \geq 1$ ) that exponentially decay is consistent with RPE coding, as seen in our data and in previous reports<sup>4,20</sup>.

**Video analysis.** During recording sessions, videos were taken at 30 frames  $s^{-1}$  of the rats as they performed the task. During the optogenetic sessions, videos were taken at 6–8 frames  $s^{-1}$ . These videos permitted analysis of movement around the behavioral chamber. We used DeepLabCut<sup>49,50</sup> in Python to determine the location of the rat's head in each frame. DeepLabCut generates a likelihood for the location of each feature in each frame, and we discarded any frames  $<0.95$ . We further processed the *x*-coordinate and *y*-coordinate traces to remove outliers above 2 s.d. of the median across moving 1-s bins. These traces were used to calculate the location of the rat within a 0.2-s window surrounding each cue onset, and the locations of the rat in 0.2-s bins from the last lick within the first 15 s after reward delivery (or 15 s even if rats were still licking) until the next cue onset for rats from the recording sessions, or from the final port exit within the first 10 s after reward delivery (or 10 s even if the rats were still in the port) until the next cue onset for rats from the optogenetic sessions. To find the average distance from the port during this time period, we found the area under the curve for distance from the port and divided by the total time. To compare this measure across sucrose and maltodextrin (or sucrose + laser) trials, we found the average across all trials of each type of rat and compared the two groups with Wilcoxon's signed-rank test. For the electrophysiology experiment, this measure was then correlated (Spearman's) to the activity of each RPE neuron in our bin of interest on each trial. To compare with shuffled data, we produced 1,000 correlations for each neuron with shuffled trial order, and compared the true mean to the distribution of means from the 1,000 shuffled populations. For the optogenetic experiment, we calculated for each rat the fractional change in distance from the port produced by the laser by dividing the difference (laser – no laser) by the no-laser value. We compared the values from these two groups with Wilcoxon's rank-sum test.

**Evolution of activity across session.** To visualize how the reward-evoked activity of neurons changed across each reward block in the blocked task, we plotted the mean activity within our bin of interest (0.75–1.95 s post-reward delivery) for five groups of three trials at a time, equally spaced throughout the completed trials of each reward (and applied the same approach to the random sucrose/maltodextrin task, as well). To assess the impact of session progress on firing rate, we pooled the activity on each trial for all neurons of interest (say, RPE cells in sessions with sucrose block first) and the proportional progress throughout the session (of total completed trials) for the respective trial, and performed a linear regression.

**Statistics and reproducibility.** Data are presented as mean  $\pm$  s.e.m. unless otherwise noted. Statistical analyses were performed in MATLAB (MathWorks) on unsmoothed data. Specific tests are noted in the text, figure legends and throughout Methods. We did not test for normality; rather, we elected to use nonparametric tests (Wilcoxon's rank-sum and signed-rank tests, all two sided) and Poisson's models (although we replicated our main result with Gaussian models on z-scored spike counts). No statistical methods were used to predetermine sample sizes, but our sample sizes are similar to those reported in previous publications<sup>8,12,22</sup>. Data collection and analysis were not performed blind to the conditions of the experiments; however, the experimenter was not in the room during data collection, and the analysis scripts were applied uniformly to all subjects. For electrophysiology experiments, there were no features on which to randomize. For optogenetic experiments, sex and experimental group were randomized across two cohorts (run consecutively on each day). Stimulus presentation was randomized and unique for each rat (except when we intentionally explored repeated stimulus presentation in the blocked sucrose/maltodextrin experiment). We excluded sessions where rats did not complete a sufficient number of trials for analysis as noted in "Optogenetic manipulations" above. For electrophysiology sessions, we included the sessions that maximized the number of neurons recorded on each wire to increase our sample size. We excluded rats if the wires, virus or optic fibers were outside the VP (or the NAc). The main finding of RPE signaling in the VP was replicated in four different tasks across two

groups of rats. The correlation between RPE signaling and task engagement was replicated in two tasks with the same group of rats. RPE signaling in the NAc was tested in only one task. The optogenetic experiments were conducted only once each.

**Reporting summary.** Further information on research design is available in the Nature Research Reporting Summary.

### Data availability

The data generated and analyzed for this manuscript are available publicly at <https://doi.org/10.12751/g-node.3lbd0c> and ref.<sup>51</sup>.

### Code availability

The code used to analyze and visualize the data in this manuscript are available as Supplementary software and online at <https://doi.org/10.12751/g-node.3lbd0c> and ref.<sup>51</sup>.

### References

48. Elber-Dorozko, L. & Loewenstein, Y. Striatal action-value neurons reconsidered. *eLife* **7**, e34248 (2018).
49. Mathis, A. et al. DeepLabCut: markerless pose estimation of user-defined body parts with deep learning. *Nat. Neurosci.* **21**, 1281–1289 (2018).
50. Nath, T. et al. Using DeepLabCut for 3D markerless pose estimation across species and behaviors. *Nat. Protoc.* **14**, 2152–2176 (2019).
51. Ottenheimer, D. J. et al. Analysis of a reward prediction error signal in ventral pallidum. *G-Node* <https://doi.org/10.12751/g-node.3lbd0c> (2020).

### Acknowledgements

This work was supported by the National Institutes of Health (grant nos. 5T32NS91018-17 (to D.J.O.), F30MH110084 (to B.A.B.), K99AA025384 (to J.M.R.), R01DA042038 and R01NS104834 (to J.Y.C.), and R01DA035943 (to P.H.J.)), by Klingenstein-Simons, MQ, NARSAD, and Whitehall (to J.Y.C.), by a NARSAD Young Investigator Award (to J.M.R.) and by the National Science Foundation Graduate Research Fellowship (grant no. DGE1746891 to D.J.O.). We thank K. Wang and X. Tong for technical assistance.

### Author contributions

D.J.O., J.M.R. and P.H.J. designed the experiments. D.J.O. collected the electrophysiology data. D.J.O., K.M.F. and T.H.K. collected the optogenetic data. B.A.B. designed and fit the models in consultation with D.J.O. D.J.O., B.A.B. and E.S. analyzed and visualized the data. D.J.O., B.A.B., J.M.R., J.Y.C. and P.H.J. interpreted the data. D.J.O., B.A.B. and P.H.J. prepared the manuscript with comments from E.S., K.M.F., T.H.K., J.M.R. and J.Y.C.

### Competing interests

The authors declare no competing interests.

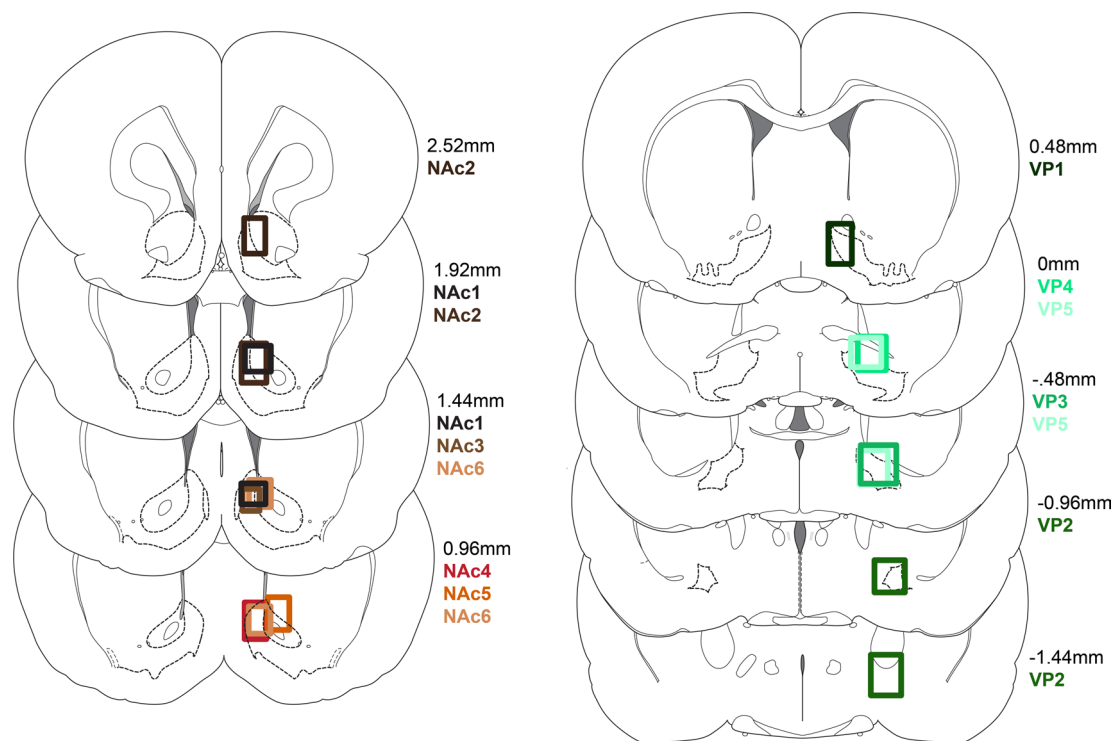
### Additional information

**Extended data** is available for this paper at <https://doi.org/10.1038/s41593-020-0688-5>.

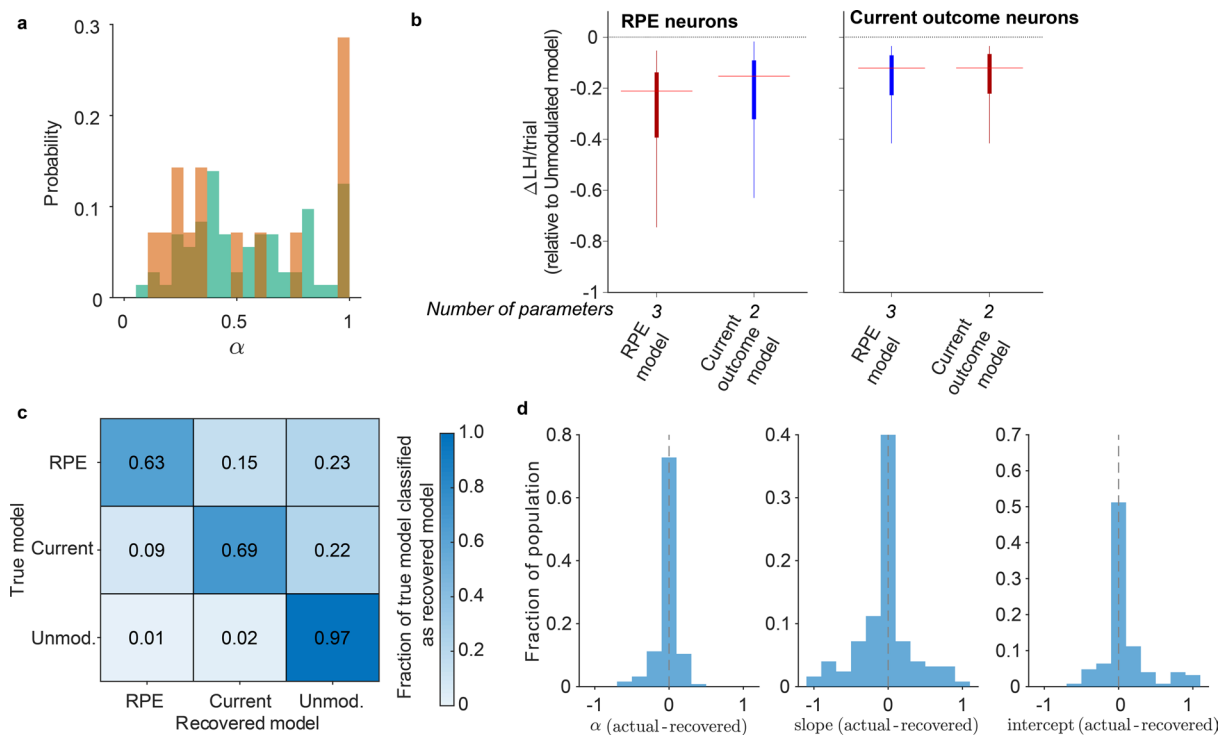
**Supplementary information** is available for this paper at <https://doi.org/10.1038/s41593-020-0688-5>.

**Correspondence and requests for materials** should be addressed to P.H.J.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

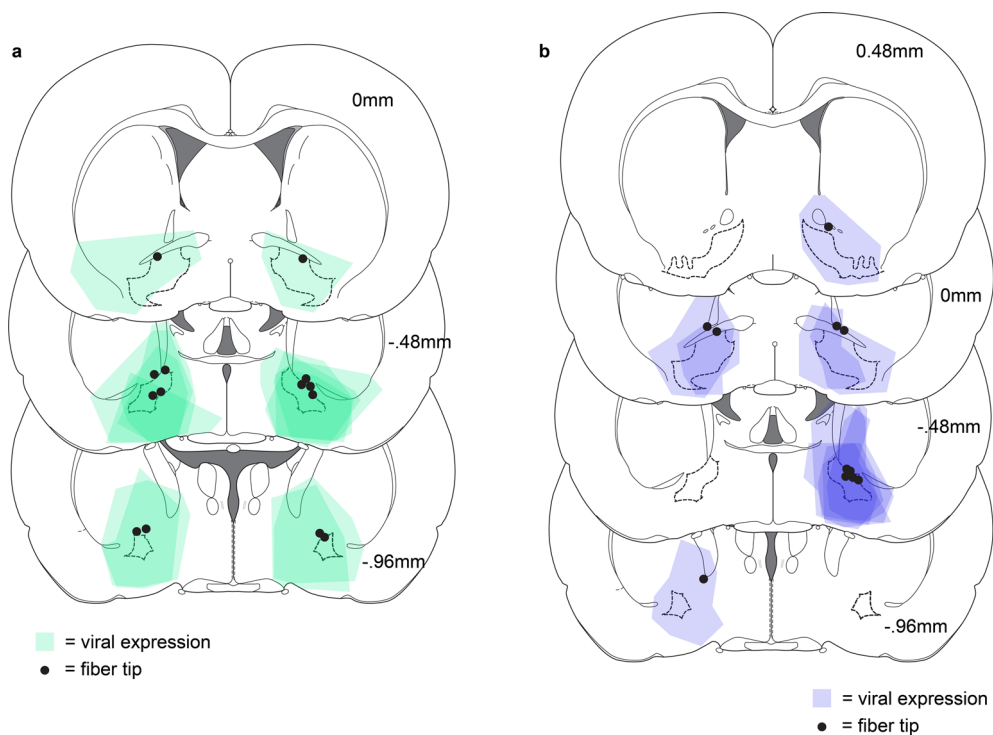


**Extended Data Fig. 1 |** Placements for random sucrose/maltodextrin, random sucrose/maltodextrin/water, and blocked sucrose/maltodextrin rats. Recording locations for nucleus accumbens (left) and ventral pallidum (right) rats.

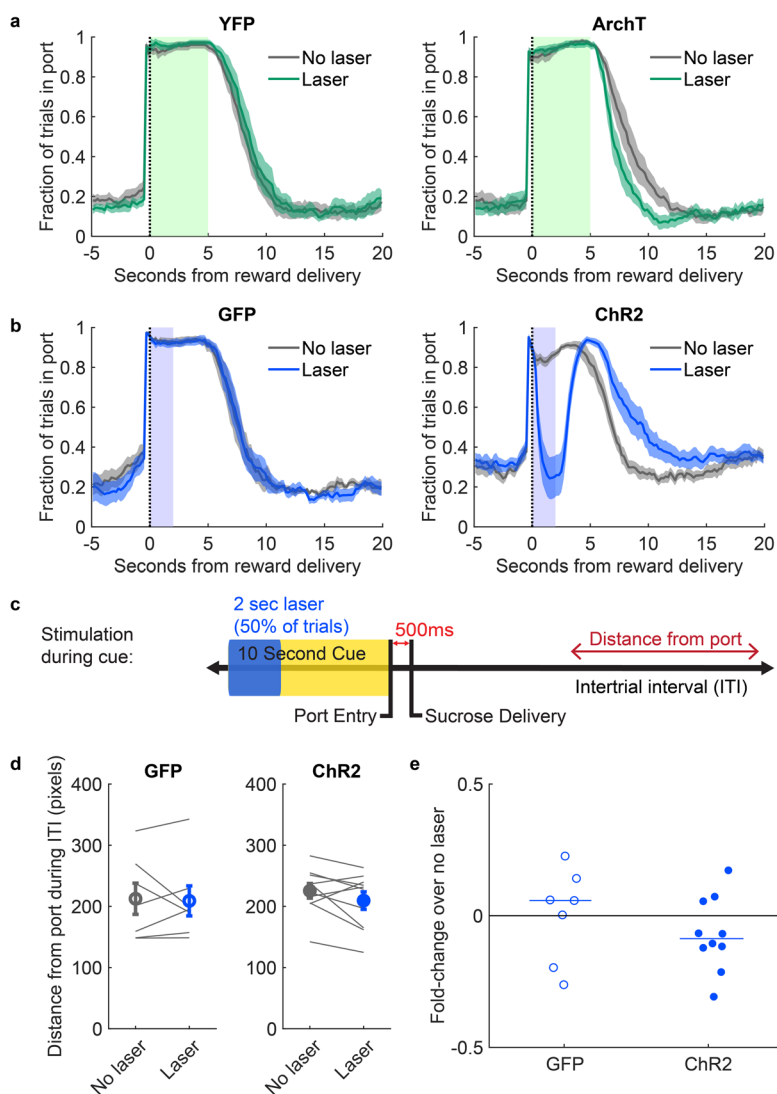


**Extended Data Fig. 2 | Evaluation of model fitting.** (a) Distribution of the learning rate,  $\alpha$ , for RPE neurons in VP (green) and NAc (orange). (b) Likelihood (LH) per trial for RPE and Current outcome neurons ( $n=72$  RPE and 126 Current outcome neurons from 5 rats) for RPE and Current outcome models, relative to the LH per trial of the Unmodulated model. Lower (more negative) indicates a better fit. Line represents median, box represents 25th and 75th percentile, and whiskers extend to 1.5 times the interquartile range. Red highlights the AIC-selected model. Median [25<sup>th</sup> to 75<sup>th</sup> percentile; min to max]  $\Delta \text{LH}/\text{trial}$  are: RPE neurons, RPE model  $-0.21$  [ $-0.39$  to  $-0.14$ ;  $-3.16$  to  $-0.05$ ], RPE neurons, Current outcome model  $-0.15$  [ $-0.32$  to  $-0.09$ ;  $-3.03$  to  $-0.02$ ], Current outcome neurons, RPE model  $-0.12$  [ $-0.23$  to  $-0.07$ ;  $-0.174$  to  $-0.03$ ], Current outcome neurons, Current outcome model  $-0.12$  [ $-0.22$  to  $-0.07$ ;  $-1.73$  to  $-0.03$ ]. Median [25<sup>th</sup>-75<sup>th</sup> percentile] LH per trial for RPE neurons was 2.29 [2.04 to 2.49] and for Current outcome neurons was 2.15 [1.92 to 2.37]. (c) Model recovery, plotted as the fraction of neurons simulated with each model recovered as that model. (d) Distribution of difference between the true value of the parameters used to simulate the neurons in (c) and the values recovered by MLE.

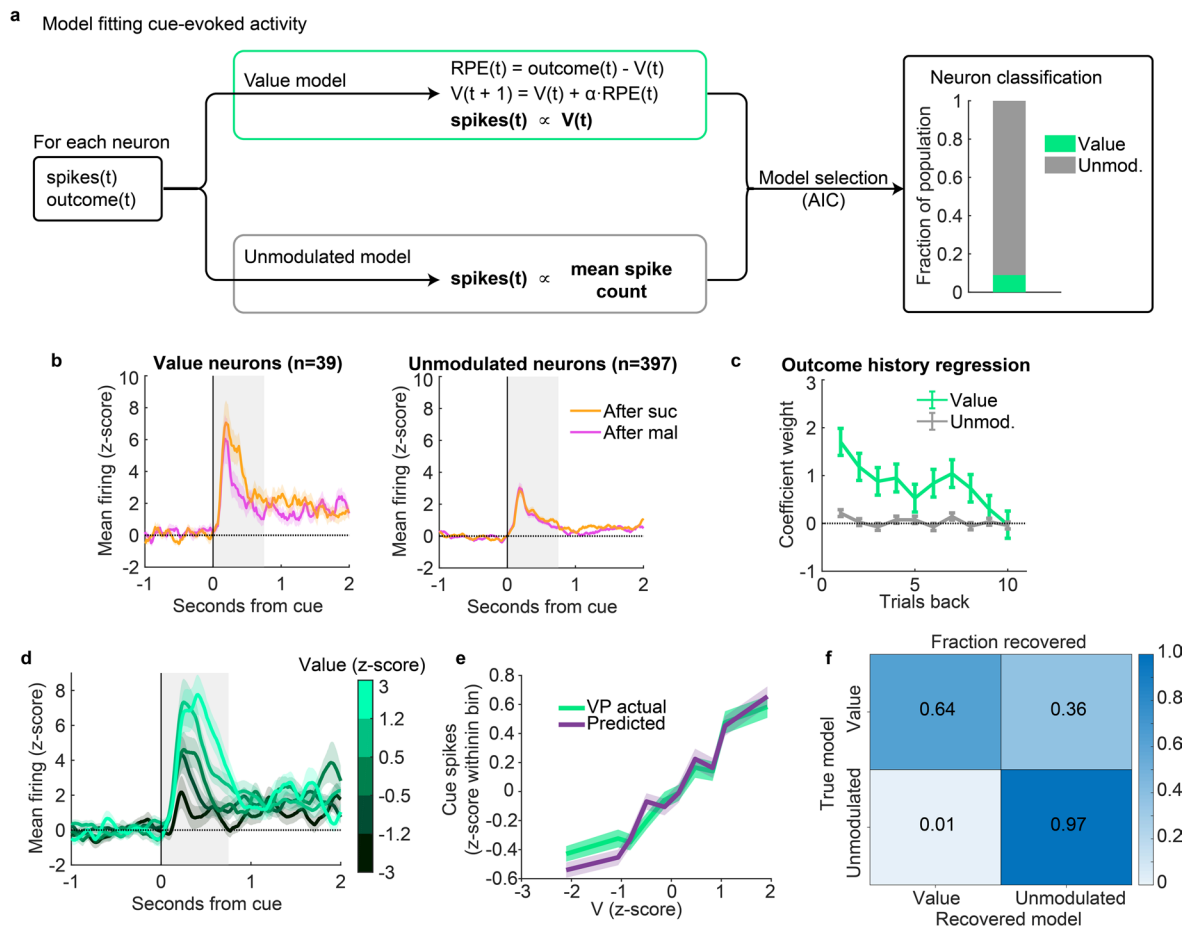




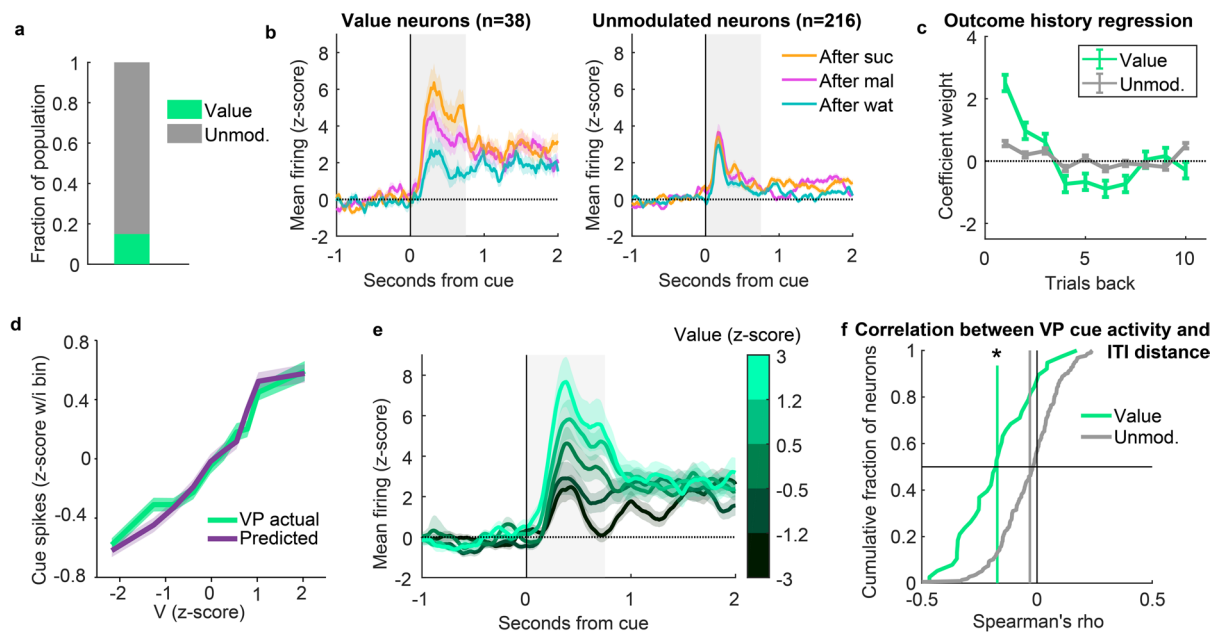
**Extended Data Fig. 3 | Placements for optogenetic experiments.** (a) Expression of ArchT3.0:YFP and fiber tip placement for the rats included in the ArchT3.0 group for the optogenetic experiment in Fig. 3. (b) Expression of ChR2:GFP and fiber tip placement for the rats included in the ChR2 group. Pattern of results remained unchanged with or without inclusion of the rat with the most caudal placement.



**Extended Data Fig. 4 | Supplemental optogenetic data.** (a) Mean( $\pm$ SEM) port occupancy in time surrounding reward delivery on laser and no laser trials for YFP (left,  $n=7$  rats) and ArchT (right,  $n=7$  rats) groups. (b) Mean( $\pm$ SEM) port occupancy in time surrounding reward delivery on laser and no laser trials for GFP (left,  $n=7$  rats) and ChR2 (right,  $n=11$  rats) groups. To account for the disruption of port occupancy by laser stimulation, we ran our distance from port analysis on the time beyond 15 s past reward delivery and found the same pattern of results. (c) Additional optogenetic experiment in ChR2 rats and controls where the 2 sec of laser stimulation was at the onset of the cue. (d) Mean( $\pm$ SEM) distance from port in the ITI following laser stimulation did not differ from no laser trials for GFP ( $p=0.94$ , Wilcoxon signed-rank test, two-sided,  $n=7$  rats) or ChR2 ( $p=0.11$ , Wilcoxon signed-rank test, two-sided,  $n=10$  rats) groups. (e) The effect of laser was similar across both groups (median: 0.06 GFP,  $n=7$  rats;  $-0.09$  ChR2,  $n=10$  rats;  $p=0.36$ , Wilcoxon rank-sum test, two-sided).

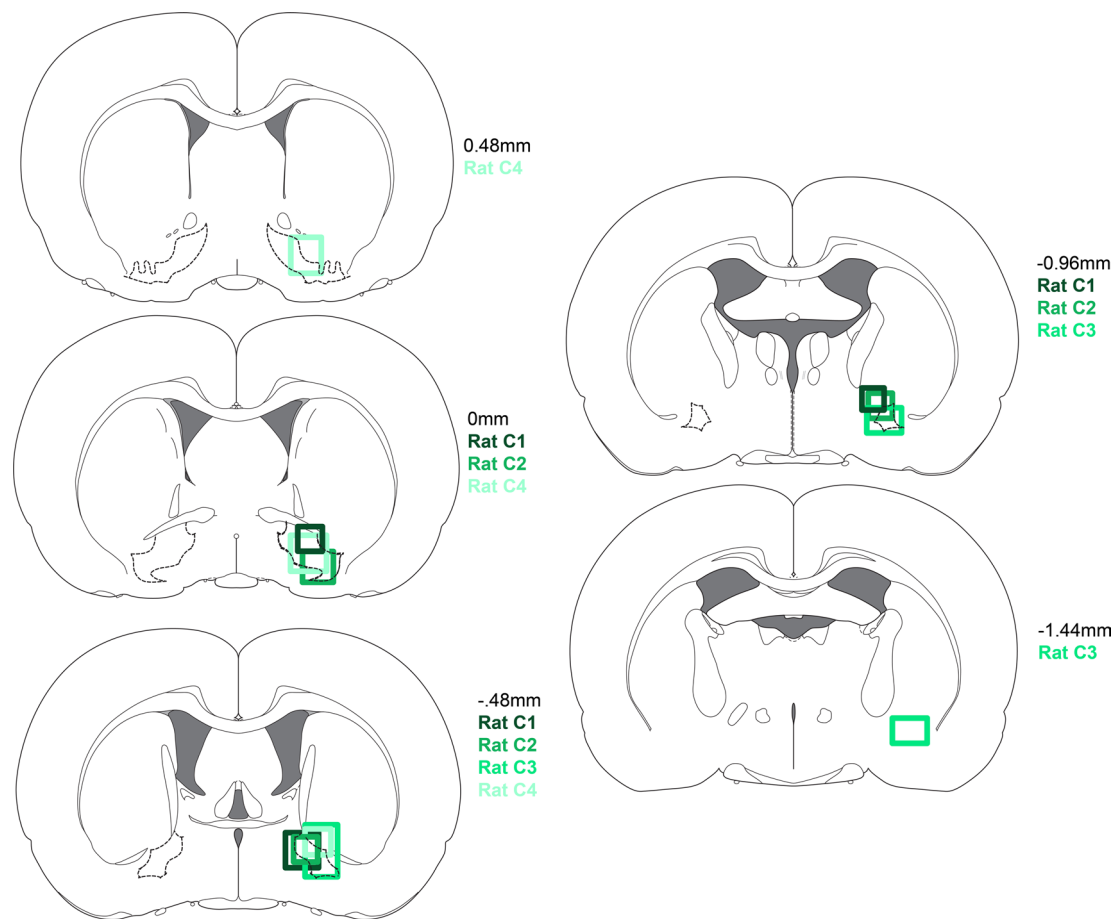


**Extended Data Fig. 5 | Value encoding in VP at the time of cue onset in the random sucrose/maltodextrin task.** (a) Schematic of model-fitting and neuron classification process. For each neuron, the reward outcome and spike count following reward delivery on each trial were used to fit two models: Value and Unmodulated. Akaike information criterion (AIC) was used to select the best model (right). (b) Mean(±SEM) activity of neurons best fit by each of the models, plotted according to previous outcome (n = 39 Value and 397 Unmodulated neurons from 5 rats). (c) Coefficients(±SE) for outcome history linear regression for each class of neurons (n = 39 Value and 397 Unmodulated neurons). (d) Mean(±SEM) activity of all Value neurons with trials binned by model-derived Value. (e) Mean(±SEM) population activity of simulated and actual Value neurons according to each trial's Value (V). (f) Model recovery, plotted as the fraction of neurons simulated with each model recovered as that model.

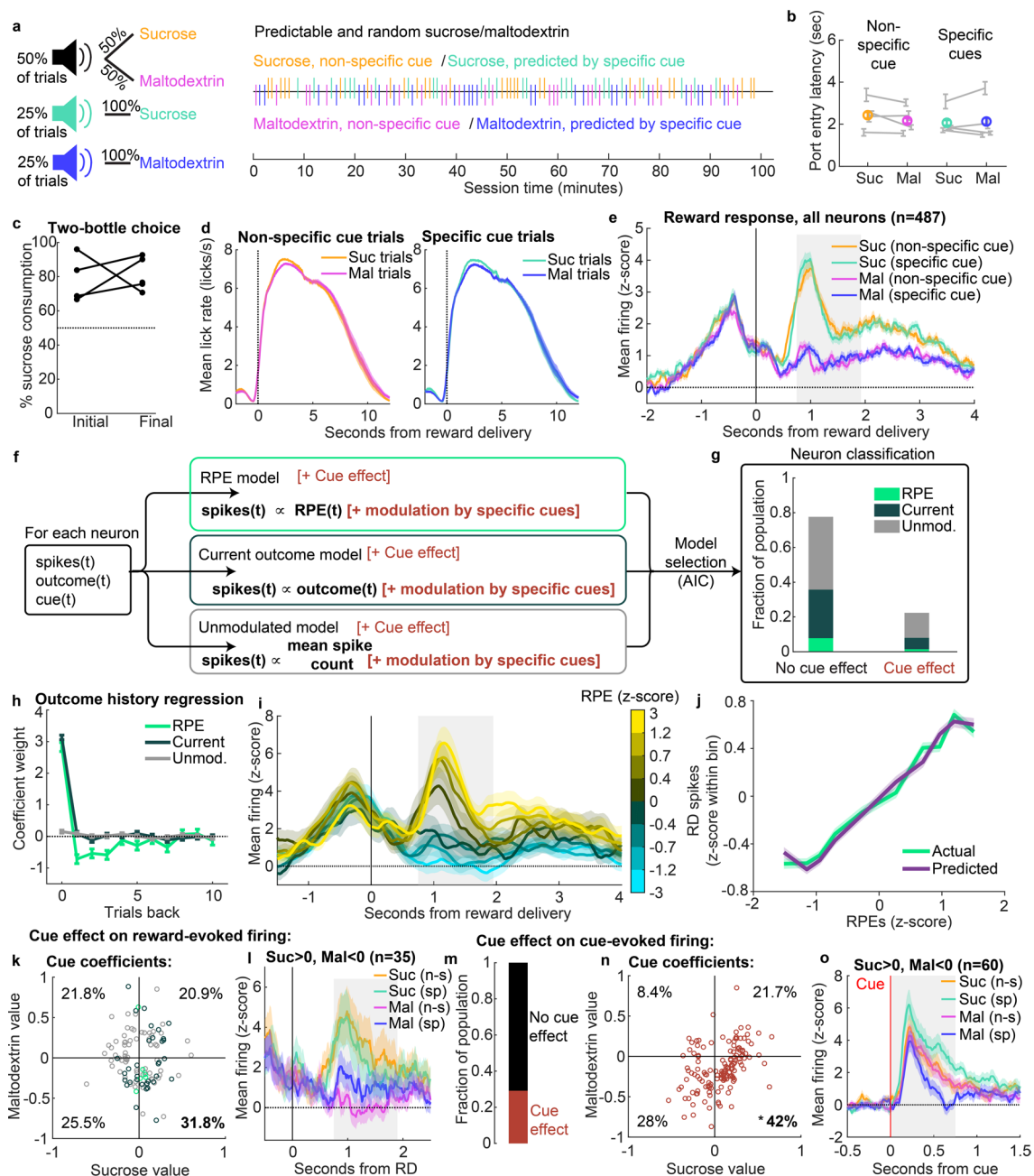


**Extended Data Fig. 6 | Value encoding at the time of cue onset in the random sucrose/maltodextrin/water task.** (a) Fraction of VP neurons best fit by the Value and Unmodulated models in the random sucrose/maltodextrin/water task. (b) Mean( $\pm$ SEM) activity of neurons best fit by each of the models, plotted according to previous outcome ( $n=38$  Value and 216 Unmodulated neurons from 3 rats). (c) Coefficients( $\pm$ SE) for outcome history linear regression for each class of neurons ( $n=38$  Value and 216 Unmodulated neurons). (d) Mean( $\pm$ SEM) population activity of simulated and actual Value neurons according to each trial's Value ( $V$ ). (e) Mean( $\pm$ SEM) activity of all Value neurons with trials binned by model-derived Value. (f) Distribution of correlations between individual VP neurons' firing rates at cue onset on each trial and the distance from the port during the previous ITI. \*  $p=0.00001$  for negative shift in mean correlation coefficient (vertical line) compared to 1000 shuffles of data for Value neurons, Wilcoxon signed-rank test, two-sided, as well as  $p=0.0000002$  for more negative coefficients for Value neurons compared to Unmodulated neurons, Wilcoxon rank-sum test, two-sided. See also Fig. 4c,d.



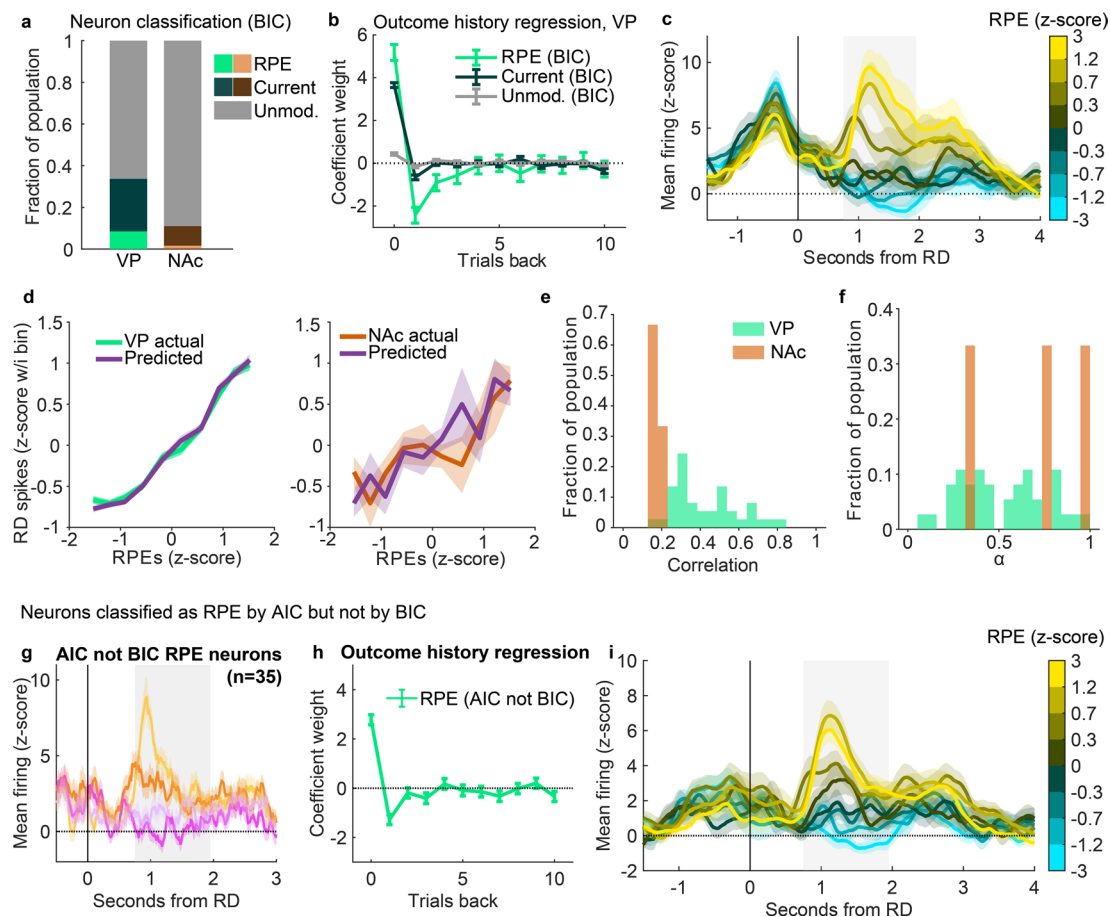


**Extended Data Fig. 7 | Placements for predictable and random sucrose/maltodextrin rats.** Recording locations for rats from predictable and random sucrose/maltodextrin experiment in Extended Data Fig. 8.



**Extended Data Fig. 8 | Impact of specific cue-derived predictions on VP firing.** (a) Task schematic: three auditory cues indicated three trial types. (b) Median latency to enter reward port following onset of cue for each trial type, plotted as the mean (+/–SEM) across all sessions for each rat (gray lines, n = 8, 9, 10, and 10 sessions for the 4 rats) and the overall mean (+/–SEM) (n = 37 sessions). (c) Percentage sucrose of total solution consumption in a two-bottle choice, before ('Initial') and after ('Final') recording (n = 4 rats). (d) Mean (+/–SEM) lick rate relative to reward delivery for each trial type (n = 37 sessions from 4 rats). (e) Mean (+/–SEM) activity of all neurons recorded in the predictable and random sucrose/maltodextrin task, aligned to reward delivery (n = 487 neurons from 4 rats). (f) Schematic of cue model-fitting. The best model (of 6 total) was selected with Akaike information criterion. (g) Fraction of the population best fit by each model. (h) Coefficients (+/–SE) for outcome history regression for each class of neurons with no cue effect (n = 38 RPE, 135 Current outcome, and 204 Unmodulated neurons). (i) Mean (+/–SEM) activity of all RPE neurons with no cue effect (n = 38 neurons). The trials for each neuron are binned according to their model-derived RPE. (j) Population activity of simulated and actual VP RPE neurons with no cue effect according to each trial's RPE value. (k) Scatterplot of each cue effect neuron's weight for specific sucrose and maltodextrin cues (n = 7 RPE, 33 Current outcome, and 70 Unmodulated cells with cue effects). The percentage of neurons falling in each quadrant is indicated. The percentage in our quadrant of interest (positive value for sucrose and negative value for maltodextrin) did not differ from chance (p = 0.1 for exact binomial test compared to null of 25%). (l) Mean (+/–SEM) activity of neurons with sucrose values > 0 and maltodextrin values < 0, consistent with a value-based cued expectation modulation. (m) Neurons with cue effects for cue-evoked signaling, rather than reward-evoked signaling, as in (g). (n) As in (k), for activity at the time of the cue rather than time of reward (n = 143 neurons with cue effects). \* = p = 0.00001 for exact binomial test compared to null of 25%. (o) As in (l), for activity at the time of the cue rather than time of reward.

## Neurons classified with BIC instead of AIC



**Extended Data Fig. 9 | Classifying neurons with BIC instead of AIC.** (a) Fraction of neurons classified as RPE, Current outcome, and Unmodulated in VP and NAc in the random sucrose/maltodextrin task using Bayesian information criterion (BIC) as the selection criterion. (b) Coefficients ( $\pm$  SE) for outcome history regression for VP neurons of each BIC subset ( $n=37$  RPE, 110 Current outcome, and 289 Unmodulated cells from 5 rats). (c) Population mean ( $\pm$  SEM) of all VP BIC RPE neurons, binned according to the model-derived RPE. (d) Mean ( $\pm$  SEM) population activity of simulated and actual BIC RPE neurons according to each trial's RPE value for VP (left) and NAc (right). (e) Distribution of correlations between model-predicted and actual spiking for all RPE neurons from each region. (f) Distribution of  $\alpha$  for RPE neurons in VP (green) and NAc (orange). (g) Mean ( $\pm$  SEM) activity of VP neurons classified as RPE by AIC but not BIC according to current and previous outcome ( $n=35$  neurons). (h) Coefficients ( $\pm$  SE) for outcome history regression for these neurons. (i) Mean ( $\pm$  SEM) activity of these neurons binned according to model-derived RPE on each trial.

## Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

### Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a	Confirmed
<input type="checkbox"/>	<input checked="" type="checkbox"/> The exact sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement
<input type="checkbox"/>	<input checked="" type="checkbox"/> A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
<input type="checkbox"/>	<input checked="" type="checkbox"/> The statistical test(s) used AND whether they are one- or two-sided <i>Only common tests should be described solely by name; describe more complex techniques in the Methods section.</i>
<input type="checkbox"/>	<input checked="" type="checkbox"/> A description of all covariates tested
<input type="checkbox"/>	<input checked="" type="checkbox"/> A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
<input type="checkbox"/>	<input checked="" type="checkbox"/> A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
<input type="checkbox"/>	<input checked="" type="checkbox"/> For null hypothesis testing, the test statistic (e.g. $F$ , $t$ , $r$ ) with confidence intervals, effect sizes, degrees of freedom and $P$ value noted <i>Give <math>P</math> values as exact values whenever suitable.</i>
<input type="checkbox"/>	<input checked="" type="checkbox"/> For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
<input checked="" type="checkbox"/>	<input type="checkbox"/> For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
<input type="checkbox"/>	<input checked="" type="checkbox"/> Estimates of effect sizes (e.g. Cohen's $d$ , Pearson's $r$ ), indicating how they were calculated

*Our web collection on [statistics for biologists](#) contains articles on many of the points above.*

### Software and code

Policy information about [availability of computer code](#)

Data collection	Custom Med Associates codes were used to direct the animal's behavior during the session and collect behavioral responses. A custom arduino-based code allowed registering of licks during the session. OmniPlex 17 (from Plexon) was used to acquire spiking information and interface with behavioral data.
Data analysis	Offline Sorter V3 from Plexon was used to sort spiking data into units. NeuroExplorer 4 (Nex Technologies) was used for event creation and initial visualization. DeepLabCut (open source) was used to analyze videos in Python 3.7. Custom scripts in MATLAB_R2020a (Mathworks) were used for all subsequent analysis and visualization of behavior and neural activity, available in Supplementary Software and at: <a href="https://doi.org/10.12751/g-node.3lbd0c">https://doi.org/10.12751/g-node.3lbd0c</a>

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

### Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

The data generated and analyzed for this manuscript (timestamps of behavioral events, task events, and individual neuron spikes) are available in Supplementary Software and publicly at: <https://doi.org/10.12751/g-node.3lbd0c>



## Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

☒ Life sciences ☐ Behavioural & social sciences ☐ Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://www.nature.com/documents/nr-reporting-summary-flat.pdf)

## Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	No statistical methods were used to pre-determine sample sizes, but our sample sizes are similar to those reported in previous publications with electrophysiology and optogenetic experiments (for example, citations 8,12,22).
Data exclusions	Rats were excluded from electrophysiology experiment analysis if no neurons were recorded from the animal or placement of the electrodes was in neither region of interest. Sessions were excluded if rats completed fewer than 15 trials of either reward (which limits the ability to evaluate event-related responding). Rats from the optogenetic experiments were excluded if one or both of the fiber optics (or viral expression) were aimed outside of ventral pallidum or if they completed fewer than 30 trials, which limited the ability to assess an impact of laser.
Replication	All experiments were repeated in multiple animals. The number and nature of reward-specific neural responses were evaluated across all animal subjects to ensure consistency. Our main finding of prediction error encoding was replicated in 4 experiments from 2 cohorts of rats. The correlation between RPE signaling and task engagement was replicated in 2 tasks with the same group of rats. RPE signaling in NAc was only tested in one task. The optogenetic experiments were only conducted once each.
Randomization	All rats were acquired at the same age and weight and single-housed, so there were no variables upon which to randomize for electrophysiology experiments. For optogenetic experiments, control and experimental groups were run simultaneously. Sex and time of day were randomized across groups.
Blinding	Data collection and analysis were not performed blind to the conditions of the experiments. For electrophysiology experiments, it was not feasible to blind the experiment to region due to the electrophysiological properties of recording from each region, or to task, which needed to be started by hand, but the experimenter was not present in the room during the experiments, minimizing the possibility for experimenter impact. All data were analyzed with the same analysis scripts, ensuring consistent treatment of different groups. For the optogenetic experiments, the experimenter was not in the room during data collection, and group assignment was randomized, minimizing the need for blinding, and, again, analysis was performed uniformly on all subjects, removing the possibility of experimenter bias.

## Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

### Materials & experimental systems

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> Antibodies
<input checked="" type="checkbox"/>	<input type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology and archaeology
<input type="checkbox"/>	<input checked="" type="checkbox"/> Animals and other organisms
<input checked="" type="checkbox"/>	<input type="checkbox"/> Human research participants
<input checked="" type="checkbox"/>	<input type="checkbox"/> Clinical data
<input checked="" type="checkbox"/>	<input type="checkbox"/> Dual use research of concern

### Methods

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging

## Animals and other organisms

Policy information about [studies involving animals](#); [ARRIVE guidelines](#) recommended for reporting animal research

Laboratory animals	Rat, Long Evans from Envigo. Age was P60 at start of experiment. Electrophysiology: all male. Optogenetics: 14 male, 17 female.
Wild animals	No wild animals were used.
Field-collected samples	No field-collected samples were used.
Ethics oversight	All experimental procedures were performed in strict accordance with protocols approved by the Animal Care and Use Committee at Johns Hopkins University.

Note that full information on the approval of the study protocol must also be provided in the manuscript.