

# Value Representations in the Orbitofrontal Cortex Drive Learning, not Choice

Kevin J. Miller<sup>1</sup>, Matthew M. Botvinick<sup>2,3,5</sup>, and Carlos D. Brody<sup>1,4,5</sup>

<sup>1</sup> Princeton Neuroscience Institute, Princeton University, Princeton, NJ, USA

<sup>2</sup> DeepMind, London, UK

<sup>3</sup> Gatsby Computational Neuroscience Unit, University College London, London, UK

<sup>4</sup> Howard Hughes Medical Institute and Department of Molecular Biology, Princeton University, Princeton NJ, USA

<sup>5</sup> Co-corresponding authors

**Humans and animals make predictions about the rewards they expect to receive in different situations. In formal models of behavior, these predictions are known as “value representations”, and they play two very different roles. Firstly, they drive *choice*: expected values of available options are compared to one another, and the best is selected. Secondly, they support *learning*: expected values are compared to rewards actually received, and future expectations are updated accordingly. A fundamental unanswered question is whether these different functions are mediated by the same or different neural mechanisms. Here we employ a recently-developed multi-step task for rats that cleanly separates learning from choosing. We address the role of the orbitofrontal cortex (OFC), a key player in value-based cognition. Electrophysiological recordings and optogenetic perturbations indicate that, contrary to prominent theories, the OFC does not directly drive choices. Instead, it supplies value information to a learning process that updates choice mechanisms elsewhere in the brain. This result places important constraints on neural architectures for learning and choosing.**

Representations of expected value play a key role in human and animal cognition (1–4). The orbitofrontal cortex (OFC) is well established as a key region for representing and using value information, but its particular role remains the subject of heated debate (5–8). One set of theories emphasizes a direct role for OFC in choice, proposing that it takes value information as input and transforms this information into a decision (9–11). Another set proposes indirect roles for OFC in choice, whether by representing choice-relevant information about the state of the world (12–14), or by providing key input to an ongoing choice process (15, 16). A final set proposes yet a more indirect role: that OFC drives learning of the representations upon which choices are based (17–19).

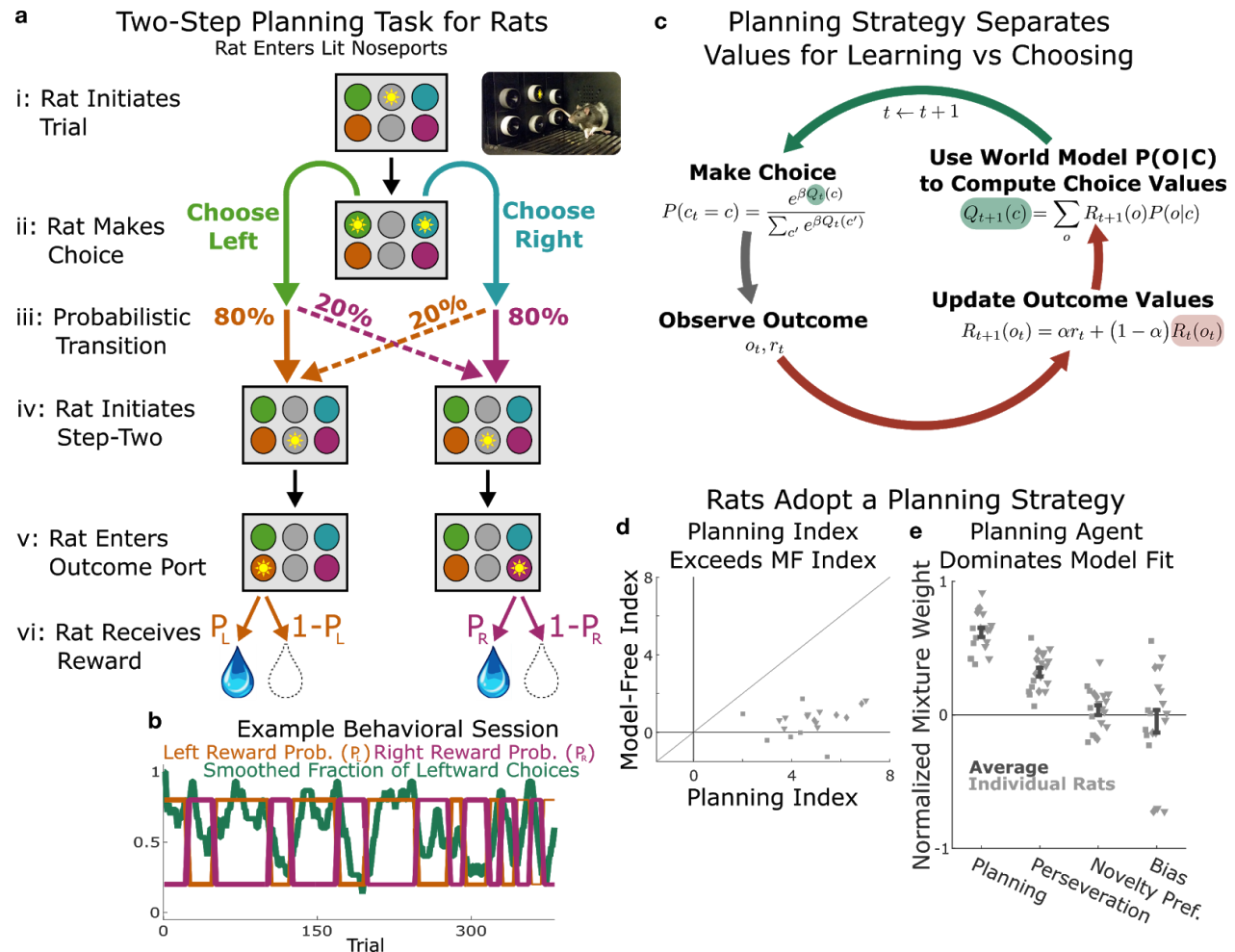
Does the OFC play a role in choice, in learning, or in both? Resolving this question has been difficult, because both ideas make similar predictions for typical laboratory tasks, in which the items to be learned about are identical to the items to be chosen between. In recent work (20), we have adapted for rats a two-stage decision task from the human literature (21) which breaks this identity. In this task, a choice made by the subject on the first step is linked probabilistically to an outcome that occurs on the second step, which in turn is linked probabilistically to reward

(**Fig. 1a**). The first-step probabilities are fixed for each rat, while the second-step probabilities are variable and must be learned. This task structure, in which choices and outcomes do not have a 1-to-1 relationship, but are linked probabilistically, combined with the particular strategy rats use to solve it, provides the critical separation between learning and choice necessary to differentiate the two. Specifically, on the first step of every trial, rats choose an option, based on values that are computed, not learned (“Use World Model” and “Make Choice” in **Fig. 1c**); while on the second step they learn the values of the outcomes (“Update Outcome Values” in **Fig. 1c**) but are led to those outcomes without making a choice.

The structure of each trial was as follows: The rat first initiated the trial by poking its nose into a neutral center port, and then selected one of two choice ports (**Fig. 1a i,ii**). One choice caused a left outcome port to become available with high probability (“common” transition), and a right outcome port to become available with low probability (“uncommon” transition), while the opposite choice reversed these probabilities (**Fig. 1a iii**). Following the initial choice, an auditory tone informed the rat which of the two outcome ports had in fact become available on that trial and, after poking into a second neutral center poke, the available outcome port was further indicated by a light (**Fig. 1a iv,v**). The rat was required to poke into the available outcome port (no choice in this step), where it received a water reward with a given probability (**Fig. 1a v,vi**). The two outcome ports differed in the probability with which they delivered reward, and these reward probabilities changed at unpredictable intervals (**Fig. 1b**). The subjects were thus required to continually update their estimates of the expected reward probability at each outcome port  $o$  (which we label  $R(o)$ ), through a learning process in which they compare their expectations to actual reward received (“Update Outcome Values” in **Fig. 1c**).

Previously (20), we have shown that rats solve the task using a particular strategy termed “model-based planning” (2, 22, 23). This strategy utilizes an internal model of action-outcome relationships to compute the expected values of the two choice ports (which we label  $Q(c)$ ). A planning strategy results in very different computational roles for outcome-port values and choice-port values. Outcome port values, ( $R$ ), are learned incrementally from recent rewards (**Fig. 1c**, “Update Outcome Values”), while choice port values ( $Q$ ) are computed based on outcome port values and the world model linking choice to outcome (which we label  $P(o|c)$ ; **Fig. 1c**, “Use World Model”), and are then used to determine the next choice (**Fig. 1c**, “Make Choice”). The values of the choice ports therefore drive choice directly, while the learned values of the outcome ports support choice only indirectly, by directly supporting learning. Consistent with previous results (20), rats in the current study adopted such a planning strategy, as indicated both by an index quantifying the extent to which the rats’ behavior is sensitive to the world model (**Fig. 1d**, planning index, Methods), and by fits of an artificial planning agent (**Fig. 1e**). This artificial agent matches rat behavior on the task, and allows us to probe the role of the OFC in two key ways. First, once the agent’s parameters have been fit to a particular rat, it can provide a trial-by-trial estimate of the value placed by that rat on the choice and on the outcome ports, which we compare below to trial-by-trial physiology data (21, 24–28). Crucially, choice port and outcome port in our task are linked only probabilistically, meaning that outcome port value and choice port value on a particular trial will be different from one another, and their neural correlates can be identified separately. Second, the agent can be altered to selectively impair information about expected values of the choice or the outcome ports, and used to

generate synthetic behavioral datasets that predict the consequences of such specific impairments. Below we compare these predictions to data from rats undergoing silencing of the OFC.

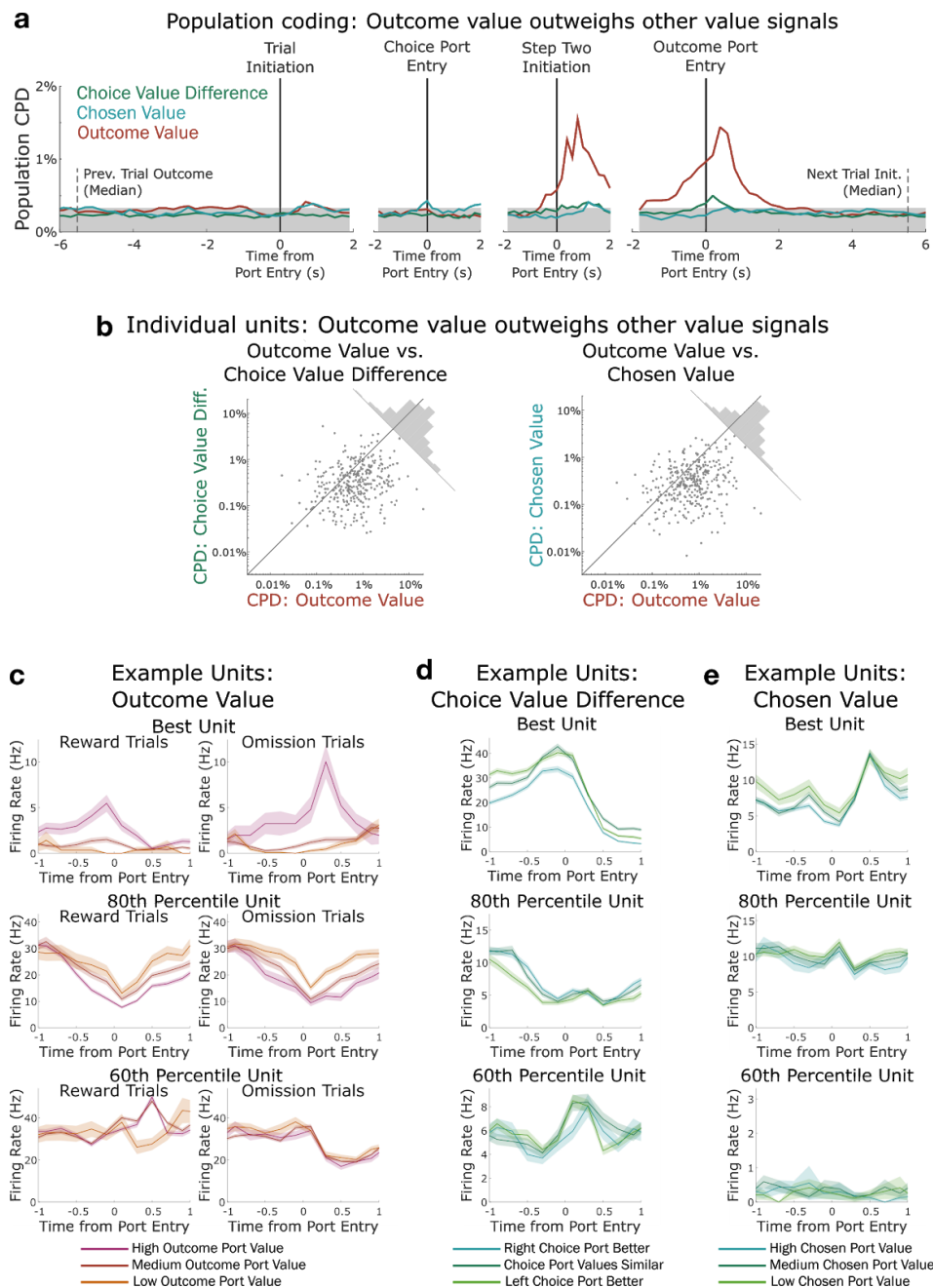


**Figure 1. Two-step task separates choice values from outcome values.** **a:** Rat two-step task. Rat initiates trial by entering the top center port (i), then chooses to enter one of two choice ports (ii). This leads to a probabilistic transition (iii) to one of two possible paths. In both paths, rat enters the bottom center port (v), causing one of two outcome ports to illuminate. Rat enters that outcome port (v), and receives a reward (vi). **b:** Example behavioral session. At unpredictable intervals, outcome port reward probabilities flip synchronously between high (80%) and low (20%). Rat adjusts choices accordingly. **c:** Schematic of the planning strategy. Agent maintains value estimates for each outcome port ( $R(o)$ ), based on a history of recent rewards at that port, agent also maintains value estimates for each choice port ( $Q(c)$ ), which are computed on each trial based on the outcome values and the world model ( $P(o|c)$ ). **d:** Planning index and model-free index, shown for electrophysiology rats ( $n=6$ , squares), optogenetics rats ( $n=9$ , triangles), and sham optogenetics rats ( $n=4$ , diamonds). These indices are calculated using a trial-history regression analysis (see Methods; ref. 20). **e:** Weights of an agent-based model fit to rats' behavioral data (see Methods; ref. 20).

Electrophysiological recordings during 51 behavioral sessions in six rats yielded 329 activity traces suitable for analysis, including both single- and multi-unit recordings (all results were similar across the two; Extended Data Fig. ED1-4). To quantify the coding properties of these units, we fit regression models to predict their spiking activity. We used regressors from

potentially learning-related events of one trial (including a model-derived estimate of outcome port value, labeled “outcome value”) as well from the choosing-related events of the subsequent trial (including two model-derived quantities: estimated value of the chosen port, and difference in estimated value between the two choice ports, labeled “chosen value” and “choice value difference”, respectively; see *Methods* for more detail and other regressors). Since the different regressors are correlated with one another (Extended Data **Fig. ED6**), we quantified the influence of each on neural firing rates using its coefficient of partial determination (CPD; also called “partial r-squared” 29, 30). CPD can be computed for a particular fit (i.e. one unit in one time bin), or for a collection of fits (aggregating variance over units, bins, or both).

First, we considered coding at the population level, computing CPD aggregated over all units for each time bin. We found that outcome-value was coded robustly, beginning at the time of entry into the neutral bottom-center port, and peaking shortly after entry into the outcome port (**Fig 2a**) at 1.5%. In contrast, population coding of choice-value-difference and chosen-value was low in all time bins, reaching a maximum of only 0.5%. Population coding of other regressors was broadly similar to previous work using other tasks (**Fig ED 4**). Next, we considered coding in individual units, aggregating data from several time bins around the time of each nose port entry. We found that a large fraction significantly modulated their firing rate according to outcome-value, with the largest fraction at outcome port entry (158/329 units, 48%; permutation test at  $p < 0.01$ ). In contrast, a relatively small fraction of units modulated their firing rate according to choice-value-difference, with the largest fraction at outcome port entry (34/329 units, 10%), or to chosen-value, with the largest fraction at choice port entry (41/329, 12%). Furthermore, the magnitude of CPD in individual units was larger for outcome-value than for the other value regressors; considering the port entry event with the strongest coding for each regressor, the mean unit had CPD for outcome-value 5.9x larger than for choice-value-difference ( $p = 10^{-16}$ , sign test; median unit 2.1x larger), and 6.9x larger than for chosen-value ( $p = 10^{-27}$ ; median unit 2.2x; **Fig. 2b**, note logarithmic axes; see also **Figs. ED2 ED3**, and **ED5**). Similar results were obtained by computing CPD for each unit over all time bins ( $p = 10^{-11}$ ,  $p = 10^{-9}$ ; **Fig. ED1**). Together, these results indicate that neural activity in OFC encodes information about values of the visited outcome port much more strongly than it encodes either type of value information about the choice ports. These signals suggest a computational role in updating expected value information at the time of outcome, rather than in selecting between choice ports.



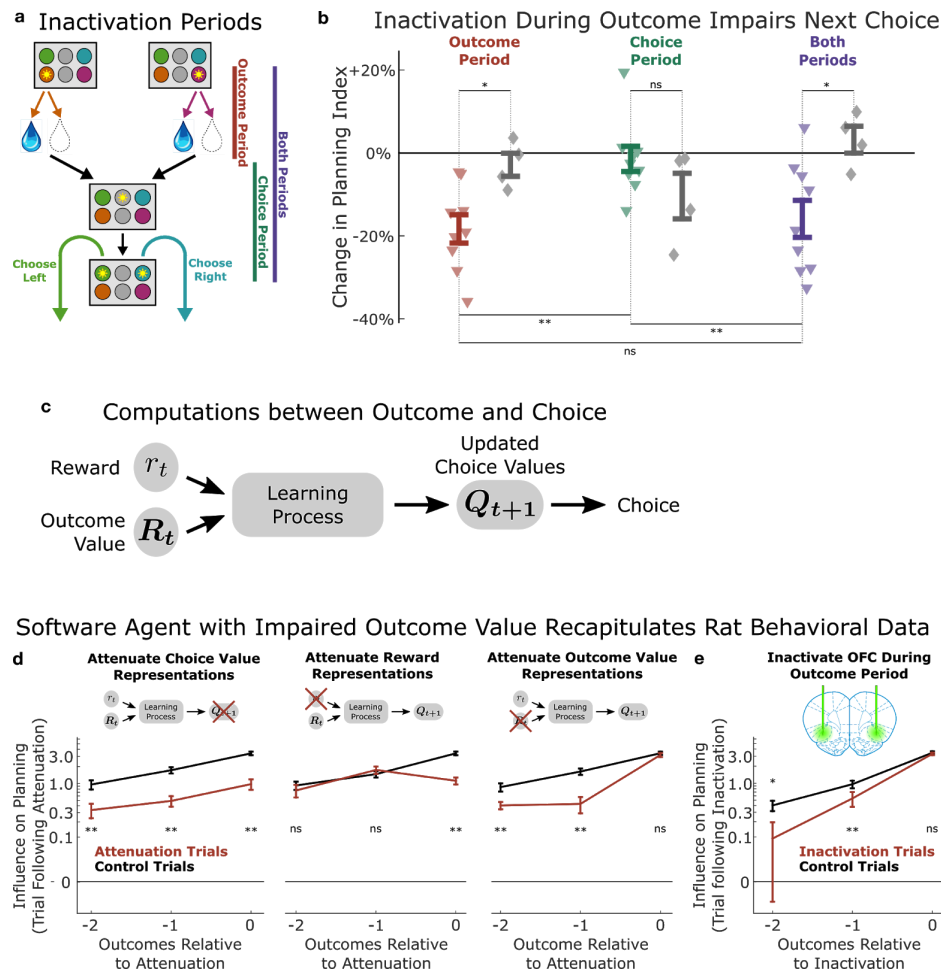
**Figure 2. OFC encodes outcome values, but not choice-related values.** **a** Timecourse of population CPD for the three value regressors. **b** Left: Scatterplot showing CPD for each unit ( $n=329$ ) for the outcome-value regressor against CPD for the choice-value-difference regressor, both computed in a one-second window centered on entry into the outcome port. Right: Scatterplot showing CPD for outcome-value (computed at outcome port entry as in panel a), against CPD for the chosen-value regressor, computed at choice port entry. **c**, Average firing rates of three example units at outcome port entry on rewarded (left) and unrewarded trials (right), separated by the expected value of the outcome port. Cells displayed are the best, 80<sup>th</sup> percentile, and 60<sup>th</sup> percentile cells (see Methods). Note differences in the activity traces by value. **d,e** Same as panel b, but cells are selected and separated based on the choice-value-difference regressor (**d**) or are separated based on the chosen-value regressor and time-locked to choice-port entry (**e**). Note relatively weak lack of separation of activity traces by value.



To help assess the causal role of the OFC's value signals, we silenced activity in the OFC using the optogenetic construct halorhodopsin (eNpHR3.0) during either the *outcome period* (beginning with entry into the outcome port and lasting until the end of reward consumption), the *choice period* (beginning at the end of reward consumption and lasting until entry into the choice port on the subsequent trial), or *both periods* (**Fig. 3a**). Previous work (20) had shown that whole-session silencing of the OFC specifically attenuates a quantity known as the “planning index” (**Fig. 1e**, see Methods), which summarizes the extent to which the rats' choices are modulated by past trials' outcomes in a way consistent with planning. Here, we found that inactivation spanning both periods similarly decreased the planning index on the subsequent trial ( $p=0.007$ , t-test; **Fig 3b**). Inactivation during the outcome period alone similarly disrupted planning ( $p=0.0006$ , **Fig 3b**), in a manner both similar to inactivation of both periods ( $p=0.47$ , paired t-test, for outcome period vs. both periods), and significantly greater than inactivation during the choice period ( $p=0.007$ , paired t-test). Choice period inactivation alone had no significant effect ( $p=0.65$ ). A control group of four rats that received sham inactivation showed no significant effect on planning index for any time period (all  $p>0.15$ ; Figure 3B, grey diamonds), and experimental and control rats differed in the effects of outcome-period and both-period inactivation ( $p=0.02$ ,  $p=0.02$ , two-sample t-tests). Together, these results confirm that silencing the OFC at the time of the outcome (and therefore the time of outcome-value coding, **Fig 2a**) is sufficient to disrupt behavior.

To assess which aspect of the behavior was affected by silencing the OFC, we perturbed our artificial planning agent in three separate ways. In the first, we attenuated choice-value representations ( $Q$ ; **Fig. 3d**, left; see Methods) on a subset of trials. In the second, we attenuated the representation of rewards actually delivered ( $r$ ; **Fig. 3d**, middle). And in the third, we attenuated outcome-value representations ( $R$ ; **Fig 3d**, right). Each of these produced a distinct pattern of behavior, most clearly visible when we separately computed the contribution to the planning index of the past three trials' rewards (**Fig. 3d**). Specifically, we found that scaling choice values on a particular trial ( $Q_{t+1}$ ) affected the influence on choice of all past rewards ( $r_t$ ,  $r_{t-1}$ , and  $r_{t-2}$ ; **Fig 3d**, left), that scaling a particular trial's reward representation ( $r_t$ ) affected only the influence of that reward (**Fig 3d**, middle), and that scaling outcome values on a particular trial ( $R_t$ ) did not affect the influence of that trial's reward ( $r_t$ ), but did affect the influence of rewards from previous trials ( $r_{t-1}$ ,  $r_{t-2}$ ; **Fig 3d** right). This is because this quantity acts as a summarized memory of the rewards of previous trials.

We performed the same analysis on the experimental data (**Fig. 3e**). We found that silencing OFC during the outcome period on a particular trial did not affect the influence of that trial's reward on the upcoming choice ( $p=0.2$ , signrank test), but that it did affect the influence of the previous two trials' outcomes ( $p=0.004$ ,  $p=0.04$ ) This pattern was consistent with the synthetic dataset in which outcome value representations had been attenuated, but not the other synthetic datasets (compare **Fig. 3e** to **Fig. 3d**). We conclude that silencing the OFC in our task predominantly impairs outcome port value information, needed for updating value expectations, but has no effect on choice port value information, needed for driving choice.



**Figure 3. Inactivation of OFC attenuates influence of outcome values.** **a:** Three time periods of inactivation. Outcome-period inactivation began when the rat entered the outcome port, and continued until the rat exited the port, or for a minimum of two seconds. Choice-period inactivation began after this outcome period, and continued until the rat entered the choice port on the next trial, or for a maximum of 15 seconds. Both-period inactivation encompassed both of these periods. **b:** Effects of inactivation on the planning index on the subsequent trial for experimental rats ( $n=9$ , colored triangles) and sham-inactivation rats ( $n=4$ , grey diamonds). Bars indicate standard errors across rats. **c:** Simplified schematic of the representations and computations that take place in our software agent between the delivery of the outcome on one trial and the choice on the next. Compare to Fig 1c. **d:** Analysis of synthetic datasets created by attenuating different representations within the software agent on a subset of trials. Each panel shows the contribution to the planning index of trial outcomes at different lags on choices, both on control trials (black) and on trials following attenuation of a representation (red). Bars indicate standard error across simulated agents (see Methods) **e:** Same analysis as in c, applied to data from optogenetic inactivation of the OFC during the outcome period.

Although learning and choosing are very different reasons to compute expected value, experiments to date have not distinguished whether value signals in the OFC drive one process, the other process, or both. The rat two-step task gave us the opportunity to separate the two roles, both in terms of coding in neural activity, and in terms of the behavioral impact of silencing that activity. Both the electrophysiological and the optogenetic results challenge the influential view that the OFC directly drives choice by representing and comparing values of available options (8–10). We find limited representation of values associated with choice options, little effect of

silencing OFC at the putative time of choice, and effects of silencing inconsistent with impairing choice values in a computational model. Instead, we find strong representation of values associated with immediately impending reward outcomes, a strong behavioral effect of silencing OFC at the time of those outcomes, and effects of that silencing that are consistent with impairing outcome values in a computational model. Our results thus support an alternative view, in which OFC activity supports choice only indirectly, by directly supporting learning (17–19). Specifically, we suggest that OFC signals expected value information to a learning process that compares this information with rewards actually received in order to update choice mechanisms. This view is also supported by existing data indicating that OFC signals expected outcomes (31, 32), and plays a role in learning (33–35), as well as the recent report that silencing OFC does not impair economic choice (36). It is also consistent with a family of proposals suggesting that OFC represents predictions about specific outcomes (37), latent state of the environment (38), or a “state space” for reinforcement learning (12–14), but adds to these proposals, suggesting that this information is used for update, rather than for choice directly. While the source and detailed nature of OFC’s representations remain an important area for future research, our results resolve an important question about their computational role. They also place a key constraint on architectures of value-based control in general: that value-based choice and value update depend on mechanisms that are at least partially separable.

### Author Contributions

KJM, MMB, and CDB conceived the project. KJM designed the experiments and analysis, with supervision from MMB and CDB. KJM carried out all experiments and analyzed the data. KJM and CDB wrote the manuscript, based on a first draft by KJM, with extensive comments from MMB.

### Acknowledgements

We thank Athena Akrami for assistance with array implant surgeries, and Jovanna Teran, Klaus Osorio, Adrian Sirko, Samantha Stein, and Lillianne Teachen for assistance with animal training. We would like to thank Nathaniel Daw and Yael Niv for helpful discussions, and Athena Akrami, Christine Constantinople, Nathaniel Daw, Cristina Domnisoru, Jeff Gauthier, Chuck Kopec, Marcelo Mattar, Bas van Opheusden, Angela Radulescu, Ben Scott, Kim Stachenfeld, and Bob Wilson for helpful comments on the manuscript.



## Methods

### Subjects

All subjects were adult male Long-Evans rats (Taconic Biosciences; Hilltop Lab Animals), placed on a restricted water schedule to motivate them to work for water rewards. Rats were housed on a reverse 12-hour light cycle and trained during the dark phase of the cycle. Rats were pair housed during behavioral training and then single housed after being implanted with microwire arrays or optical fiber implants. Animal use procedures were approved by the Princeton University Institutional Animal Care and Use Committee and carried out in accordance with NIH standards.

### Two-Step Behavioral Task

Rats were trained on a two-step behavioral task, following a shaping procedure which has been previously described(20). Rats performed the task in custom behavioral chambers containing six “nose ports” arranged in two rows of three, each outfitted with a white LED for delivering visual stimuli, as well as an infrared LED and phototransistor for detecting rats’ entries into the port. The left and right ports in the bottom row also contained sipper tubes for delivering water rewards. The rat initiated each trial by entering the illuminated top center port, causing the two top side ports (“choice ports”) to illuminate. The rat then made his choice by entering one of these ports. Immediately upon entry into a choice port, two things happened: the bottom center port light illuminated, and one of two possible sounds began to play, indicating which of the two bottom side ports (“outcome ports”) would eventually be illuminated. The rat then entered the bottom center port, which caused the appropriate outcome port to illuminate. Finally, the rat entered the outcome port which illuminated, and received either a water reward or an omission. Once the rat had consumed the reward, a trial-end sound played, and the top center port illuminated again to indicate that the next trial was ready.

The selection of each choice port led to one of the outcome ports becoming available with 80% probability (common transition), and to the other becoming available with 20% probability (uncommon transition). These probabilities were counterbalanced across rats, but kept fixed within rat for the entirety of the animal’s experience with the task. The probability that entry into each bottom side port would result in reward switched in blocks. In each block one port resulted in reward 80% of the time, and the other port resulted in reward 20% of the time. Block shifts happened unpredictably, with a minimum block length of 10 trials and a 2% probability of block change on each subsequent trial.

### Analysis of Behavioral Data: Planning Index & Model-Free Index

We quantify the effect of past trials and their outcomes on future decisions using a logistic regression analysis based on previous trials and their outcomes(20, 39). We define vectors for each of the four possible trial outcomes: common-reward (CR), common-omission (CO), uncommon-reward (UR), and uncommon-omission (UO), each taking on a value of +1 for trials of their type where the rat selected the left choice port, a value of -1 for trials of their type where the rat selected the right choice port, and a value of 0 for trials of other types. We define the following regression model:

$$\log \left( \frac{P_{left}(t)}{P_{right}(t)} \right) = \sum_{\tau=1}^T \beta_{CR}(\tau) \cdot CR(t - \tau) + \sum_{\tau=1}^T \beta_{CO}(\tau) \cdot CO(t - \tau) + \sum_{\tau=1}^T \beta_{UR}(\tau) \cdot UR(t - \tau) + \sum_{\tau=1}^T \beta_{UO}(\tau) \cdot UO(t - \tau) \quad (1)$$

where  $\beta_{cr}$ ,  $\beta_{co}$ ,  $\beta_{ur}$ , and  $\beta_{uo}$  are vectors of regression weights which quantify the tendency to repeat on the next trial a choice that was made  $\tau$  trials ago and resulted in the outcome of their type, and  $T$  is a hyperparameter governing the number of past trials used by the model to predict upcoming choice, which was set to 3 for all analyses.

We expect model-free agents to repeat choices which lead to reward and switch away from those which lead to omissions(21), so we define a model-free index for a dataset as the sum of the appropriate weights from a regression model fit to that dataset:

$$ModelFreeIndex = \sum_{\tau=1}^T [\beta_{CR}(\tau) + \beta_{UR}(\tau)] - \sum_{\tau=1}^T [\beta_{UO}(\tau) + \beta_{CO}(\tau)] \quad (2)$$

We expect that planning agents will show the opposite pattern after uncommon transition trials, since the uncommon transition from one choice is the common transition from the other choice. We define a planning index:

$$PlanningIndex = \sum_{\tau=1}^T [\beta_{CR}(\tau) - \beta_{UR}(\tau)] + \sum_{\tau=1}^T [\beta_{UO}(\tau) - \beta_{CO}(\tau)] \quad (3)$$

We also compute the main effect of past choices on future choice:

$$PastChoices = \sum_{\tau=1}^T [\beta_{CR}(\tau) + \beta_{UR}(\tau) + \beta_{UO}(\tau) + \beta_{CO}(\tau)] \quad (4)$$

### Behavior Model

We model behavior and obtain trial-by-trial estimates of value signals using an agent-based computational model which we have previously shown to provide a good explanation of rat behavior on the two-step task(20). This model adopts the mixture-of-agents approach, in which each rat's behavior is described as resulting from the influence of a weighted average of several different "agents" implementing different behavioral strategies to solve the task. On each trial, each agent  $A$  computes a value,  $Q_A(c)$ , for each of the two available choices  $c$ , and the combined model makes a decision according to a weighted average of the various strategies' values,  $Q_{total}(c)$ :

$$Q_{total}(c) = \sum_{A \in \{agents\}} \beta_A Q_A(c) \quad (5)$$

$$\pi(c) = \frac{e^{Q_{total}(c)}}{\sum_{c'} e^{Q_{total}(c')}} \quad (5)$$

where the  $\beta$ 's are weighting parameters determining the influence of each agent, and  $\pi(c)$  is the probability that the mixture-of-agents will select choice  $c$  on that trial. The model which we have previously shown to provide the best explanation of rat's behavior contains four such agents: model-based temporal difference learning, novelty preference, perseveration, and bias.

**Model-Based Temporal Difference Learning.** Model-based temporal difference learning is a planning strategy, which maintains separate estimates of the probability with which each action (selecting the left or the right choice port) will lead to each outcome (the left or the right outcome port becoming available),  $P(o|a)$ , as well as the probability,  $R(o)$ , with which each outcome will lead to reward. This strategy assigns values to the actions by combining these probabilities to compute the expected probability with which selection of each action will ultimately lead to reward:

$$Q_{plan}(c) = \sum_o R(o)P(o|c) \quad (6)$$

At the beginning of each session, the reward estimate  $R(o)$  is initialized to 0.5 for both outcomes, and the transition estimate  $P(o|c)$  is set to the true transition function for the rat being modeled (0.8 for common and 0.2 for uncommon transitions). After each trial, the reward estimate for both outcomes is updated according to

$$R(o) \leftarrow \begin{cases} (1 - \alpha) \cdot R(o) + \alpha \cdot r_t, & \text{for } o = o_t \\ (1 - \alpha) \cdot R(o) - \alpha \cdot r_t, & \text{for } o \neq o_t \end{cases} \quad (7)$$

where  $o_t$  is the outcome that was observed on that trial,  $r_t$  is a binary variable indicating reward delivery, and  $\alpha$  is a learning rate parameter constrained to lie between zero and one.

**Novelty Preference.** The novelty preference agent follows an “uncommon-stay/common switch” pattern, which tends to repeat choices when they lead to uncommon transitions on the previous trial, and to switch away from them when they lead to common transitions. Note that some rats have positive values of the  $\beta_{np}$  parameter weighting this agent (novelty preferring) while others have negative values (novelty averse; see **Fig 1e**):

$$Q_{np}(c_t) \leftarrow \begin{cases} 0, & \text{common transition trials} \\ 1, & \text{uncommon transition trials} \end{cases} \quad (8)$$

$$Q_{np}(c \neq c_t) \leftarrow 1 - Q_{np}(c_t)$$

**Perseveration.** Perseveration is a pattern which tends to repeat the choice that was made on the previous trial, regardless of whether it led to a common or an uncommon transition, and regardless of whether or not it led to reward.

$$Q_{persev}(c_t) \leftarrow 1$$

$$Q_{persev}(c \neq c_t) \leftarrow 0 \quad (9)$$

**Bias.** Bias is a pattern which tends to select the same choice port on every trial. Its value function is therefore static, with the extent and direction of the bias being governed by the magnitude and sign of this strategy’s weighting parameter  $\beta_{bias}$ .

$$Q_{bias}(left) = 1$$

$$Q_{bias}(right) = -1 \quad (10)$$

### Model Fitting

We implemented the model described above using the probabilistic programming language Stan(40, 41), and performed maximum-a-posteriori fits using weakly informative priors on all parameters(42). The prior over the weighting parameters  $\beta$  was normal with mean 0 and sd 0.5, and the prior over  $\alpha$  was a beta distribution with  $a=b=3$ . For ease of comparison, we normalize the weighting parameters  $\beta_{plan}$ ,  $\beta_{np}$ , and  $\beta_{persev}$ , dividing each by the standard deviation of its agent’s associated values ( $Q_{plan}$ ,  $Q_{np}$ , and  $Q_{persev}$ ) taken across trials. Since each weighting parameter affects behavior only by scaling the value output by its agent, this technique brings the weights into a common scale and facilitates interpretation of their relative magnitudes, analogous to the use of standardized coefficients in regression models.

### Surgery: Microwire Array Implants

Six rats were implanted with microwire arrays (Tucker-David Technologies) targeting OFC unilaterally. Arrays contained tungsten microwires 4.5mm long and 50  $\mu$ m in diameter, cut at a 60° angle at the tips. Wires were arranged in four rows of eight, with spacing 250  $\mu$ m within-row and 375  $\mu$ m between rows, for a total of 32 wires in a 1.125 mm by 1.75 mm rectangle. Target coordinates for the implant with respect to bregma were 3.1-4.2mm anterior, 2.4-4.2mm lateral, and 5.2mm ventral (~4.2mm ventral to brain surface at the posterior-middle of the array).

In order to expose enough of the skull for a craniotomy in this location, the jaw muscle was carefully resected from the lateral skull ridge in the area near the target coordinates. Dimpling of the brain surface was minimized following

procedures described in more detail elsewhere(43). Briefly, a bolus of petroleum jelly (Puralube, Dechra Veterinary Products) was placed in the center of the craniotomy to protect it, while cyanoacrylate glue (Vetbond, 3M) was used to adhere the pia mater to the skull at the periphery. The petroleum jelly was then removed, and the microwire array inserted slowly into the brain. Rats recovered for a minimum of one week, with ad lib access to food and water, before returning to training.

### *Electrophysiological Recordings*

Once rats had recovered from surgery, recording sessions were performed in a behavioral chamber outfitted with a 32 channel recording system (Neuralynx). Spiking data was acquired using a bandpass filter between 600 and 6000 Hz and a spike detection threshold of 30  $\mu$ V. Clusters were manually cut (Spikesort 3D, Neuralynx), and both single- and multi-units were considered, on the condition that they showed an average firing rate greater than 1 Hz.

### *Analysis of Electrophysiology Data*

To determine the extent to which different variables were encoded in the neural signal, we fit a series of regression models to our spiking data. Models were fit to the spike counts emitted by each unit in 200 ms time bins taken relative to the four noseport entry events that made up each trial. There were eight total regressors, defined relative to pairs of adjacent trials, and consisting of the choice port selected (left or right), the outcome port visited (left or right), the reward received (reward or omission), the interaction between outcome port and reward, and the expected value of the outcome port visited ( $V$ ) for the first trial, and the choice port selected, the value difference between the choice ports ( $Q(left) - Q(right)$ ), and the value of the choice port selected ( $Q(chosen)$ ) for the subsequent trial. These last three regressors were obtained using the agent-based computational model described above, with parameters fit separately to each rat's behavioral data. Regressors were z-scored to facilitate comparison of fit regression weights. Models were fit using the Matlab function `lassoglm` using a Poisson noise model and L1 regularization parameters (pure lasso regression;  $\alpha = 1$ ,  $\lambda = 10^{-4}$ ) sufficient to yield a non-null model for all units.

In our task, many of these regressors were correlated with one another (**Fig. ED6**), so we quantify encoding using the coefficient of partial determination (CPD; also known as partial r-squared) associated with each(29, 30). This measure quantifies the fraction of variance explained by each regressor, once the variance explained by all other regressors has been taken account of:

$$CPD(X_i, u, t) = \frac{(SSE(X_{-i}, u, t) - SSE(X_{all}, u, t))}{SSE(X_{-i}, u, t)} \quad (11)$$

where  $u$  refers to a particular unit,  $t$  refers to a particular time bin, and  $SSE(X_{all})$  refers to the sum-squared-error of a regression model considering all eight regressors described above, and  $SSE(X_{-i})$  refers to the sum-squared-error of a model considering the seven regressors other than  $X_i$ . We compute total CPD for each unit by summing the SSE associated with the regression models for that unit for all time bins:

$$CPD(X_i, u) = \frac{\sum_t (SSE(X_{-i}, u, t) - SSE(X_{all}, u, t))}{\sum_t SSE(X_{-i}, u, t)} \quad (12)$$

We report this measure both for each unit taken over all time bins (Fig ED1), as well as for the case where the sum is taken over the five bins making up a 1s time window centered on a particular port entry event (top neutral center port, choice port, bottom neutral center port, or outcome port). Cells were labeled as significantly encoding a regressor if the CPD for that regressor exceeded that of more than 99% of CPDs computed based on datasets with circularly permuted trial labels. We use permuted, rather than shuffled, labels in order to preserve trial-by-trial correlational structure. Differences in the strength of encoding of different regressors were assessed using a sign test on the differences in CPD.

We compute total population CPD for a particular time bin in an analogous way:

$$CPD(X_i, t) = \frac{\sum_u (SSE(X_{-i}, u, t) - SSE(X_{all}, u, t))}{\sum_u SSE(X_{-i}, u, t)} \quad (13)$$

Time bins were labeled as significantly encoding a regressor if the population CPD for that time bin exceeded the largest population CPD for any regressor in more than 99% of the datasets with permuted trial labels. We compute total CPD for each regressor by summing SSE over both units and time bins, and assess the significance of differences in this quantity by comparing them to differences in total CPD computed for the shuffled datasets.

### ***Surgery: Optical Fiber Implant and Virus Injection***

Rats were implanted with sharpened fiber optics and received virus injections following procedures similar to those described previously(43–45), and documented in detail on the Brody lab website ([http://brodywiki.princeton.edu/wiki/index.php/Etching\\_Fiber\\_Optics](http://brodywiki.princeton.edu/wiki/index.php/Etching_Fiber_Optics)). A 50/125  $\mu\text{m}$  LC-LC duplex fiber cable (Fiber Cables) was dissected to produce four blunt fiber segments with LC connectors. These segments were then sharpened by immersing them in hydrofluoric acid and slowly retracting them using a custom-built motorized jig attached to a micromanipulator (Narashige International) holding the fiber. Each rat was implanted with two sharpened fibers, in order to target OFC bilaterally. Target coordinates with respect to bregma were 3.5mm anterior, 2.5mm lateral, 5mm ventral. Fibers were angled 10 degrees laterally, to make space for female-female LC connectors which were attached to each and included as part of the implant.

Four rats were implanted with sharpened optical fibers only, but received no injection of virus. These rats served as uninfected controls.

Nine additional rats received both fiber implants as well as injections of a virus (AAV5-CaMKII $\alpha$ -eNpHR3.0-eYFP; UNC Vector Core) into the OFC to drive expression of the light-activated inhibitory opsin eNpHR3.0. Virus was loaded into a glass micropipette mounted into a Nanoject III (Drummond Scientific), which was used for injections. Injections involved five tracks arranged in a plus-shape, with spacing 500 $\mu\text{m}$ . The center track was located 3.5mm anterior and 2.5mm lateral to bregma, and all tracks extended from 4.3 to 5.7mm ventral to bregma. In each track, 15 injections of 23 nL were made at 100 $\mu\text{m}$  intervals, pausing for ten seconds between injections, and for one minute at the bottom of each track. In total 1.7  $\mu\text{L}$  of virus were delivered to each hemisphere over a period of about 20 minutes.

Rats recovered for a minimum of one week, with ad lib access to food and water, before returning to training. Rats with virus injections returned to training, but did not begin inactivation experiments until a minimum of six weeks had passed, to allow for virus expression.

### ***Optogenetic Perturbation Experiments***

During inactivation experiments, rats performed the task in a behavioral chamber outfitted with a dual fiber optic patch cable connected to a splitter and a single-fiber commutator (Princeton) mounted in the ceiling. This fiber was coupled to a 200 mW 532 nm laser (OEM Laser Systems) under the control of a mechanical shutter (ThorLabs) by way of a fiber port (ThorLabs). The laser power was tuned such that each of the two fibers entering the implant received between 25 and 30 mW of light when the shutter was open.

Each rat received several sessions in which the shutter remained closed, in order to acclimate to performing the task while tethered. Once the rat showed behavioral performance while tethered that was similar to his performance before the implant, inactivation sessions began. During these sessions, the laser shutter was opened (causing light to flow into the implant, activate eNpHR3.0 and silence neural activity) on 7% of trials each in one of three time periods. “Outcome period” inactivation began when the rat entered the bottom center port at the end of the trial, and ended either when the rat had left the port and remained out for a minimum of 500 ms, or after 2.5 s. “Choice period” inactivation began at the end of the outcome period and lasted until the rat entered the choice port on the following trial. “Both period” inactivation encompassed both the outcome period and the choice period. The total duration of the inactivation therefore depended in part on the movement times of the rat, and was somewhat variable from trials to trial (**Fig ED8**). If a scheduled inactivation would last more than 15 s, inactivation was terminated, and that trial was excluded from analysis. Due to constraints of the bControl software, inactivation was only performed on even-numbered trials.



## Analysis of Optogenetic Effects on Behavior

We quantify the effects of optogenetic inhibition on behavior by computing separately the planning index for trials following inactivation of each type (outcome period, choice period, both periods) and for control trials. Specifically, we fit the trial history regression model of Equation 1 with a separate set of weights for trials following inactivation of each type:

$$\log \left( \frac{P_{left}(t)}{P_{right}(t)} \right) = \sum_{\tau=1}^T \beta_{CR,i}(\tau) \cdot CR(t-\tau) + \sum_{\tau=1}^T \beta_{CO,i}(\tau) \cdot CO(t-\tau) + \sum_{\tau=1}^T \beta_{UR,i}(\tau) \cdot UR(t-\tau) + \sum_{\tau=1}^T \beta_{UO,i}(\tau) \cdot UO(t-\tau) \quad (14)$$

$$i = \begin{cases} ctrl, & \text{trial } t-1 \text{ was control trial} \\ out, & \text{trial } t-1 \text{ had outcome period inactivation} \\ ch, & \text{trial } t-1 \text{ had choice period inactivation} \\ both, & \text{trial } t-1 \text{ had both} \end{cases} \quad (15)$$

We used maximum a posteriori fitting in which the priors were Normal(0,1) for weights corresponding to control trials, and Normal( $\beta_{X,ctrl}$ , 1) for weights corresponding to inactivation trials, where  $\beta_{X,ctrl}$  is the corresponding control trial weight – e.g. the prior for  $\beta_{CR,out}$  is Normal( $\beta_{CR,ctrl}$ (1), 1). This prior embodies the belief that inactivation is most likely to have no effect on behavior, and that any effect is equally likely to be positive or negative with respect to each  $\beta$ . This ensures that our priors cannot induce any spurious differences between control and inactivation conditions into the parameter estimates. We then compute a planning index separately for the weights of each type, modifying equation 3:

$$PlanningIndex_i = \sum_{\tau=1}^T [\beta_{CR,i}(\tau) - \beta_{UR,i}(\tau)] + \sum_{\tau=1}^T [\beta_{UO,i}(\tau) - \beta_{CO,i}(\tau)] \quad (16)$$

We compute the relative change in planning index for each inactivation condition:  $(PlanningIndex_i - PlanningIndex_{ctrl}) / PlanningIndex_{ctrl}$  and report three types of significance tests on this quantity. First, we test for each inactivation condition the hypothesis that there was a significant change in the planning index, reporting the results of a one-sample t-test over rats. Next, we test the hypothesis that different inactivation conditions had effects of different sizes on the planning index, reporting a paired t-test over rats. Finally, we test the hypothesis for each condition that inactivation had a different effect than sham inactivation (conducted in rats which had not received virus injections to deliver eNpHR3.0), reporting a two-sample t-test.

To test the hypothesis that inactivation specifically impairs the effect of distant past outcomes on upcoming choice, we break down the planning index for each condition by the index of the weights the contribute to it:

$$PlanningIndex_i(\tau) = [\beta_{CR,i}(\tau) - \beta_{UR,i}(\tau)] + [\beta_{UO,i}(\tau) - \beta_{CO,i}(\tau)] \quad (17)$$

We report these trial-lagged planning indices for each inactivation condition, and assess the significance of the difference between inactivation and control conditions at each lag using a signrank test across rats.

## Synthetic Datasets

To generate synthetic datasets for comparison to optogenetic inactivation data, we generalized the behavioral model to separate the contributions of representations of expected value and of immediate reward. In particular, we replaced the learning equation within the model-based RL agent (Equation 7) with the following:

$$V(o) \leftarrow \begin{cases} \alpha_{value} \cdot V(o) + \alpha_{reward} \cdot R_t + (1 - \alpha_{value} - \alpha_{reward}) \cdot E[V], & \text{for } o = o_t \\ \alpha_{value} \cdot V(o) - \alpha_{reward} \cdot R_t + (1 - \alpha_{value} - \alpha_{reward}) \cdot E[V], & \text{for } o \neq o_t \end{cases} \quad (18)$$

where  $\alpha_{value}$  and  $\alpha_{reward}$  are separate learning rate parameters, constrained to be nonnegative and to have a sum no larger than one, and  $E[V]$  represents the expected reward of a random-choice policy on the task, which in the case of our task is equal to 0.5.

To generate synthetic datasets in which silencing the OFC impairs choice value representations, outcome value representations, or reward representations, we decrease the parameter  $\beta_{plan}$ ,  $\alpha_{value}$ , or  $\alpha_{reward}$ , respectively. Specifically, we first fit the model to the dataset for each rat in the optogenetics experiment (n=9) as above (i.e. using equation 7 as the learning rule) to obtain maximum a posteriori parameters. We translated these parameters to the optogenetics (equation 18) version of the model by setting  $\alpha_{value}$  equal to the fit parameter  $\alpha$  and  $\alpha_{reward}$  equal to  $1 - \alpha$ . We then generated four synthetic datasets for each rat. For the control dataset, the fit parameters were used on trials of all types, regardless of whether inhibition of OFC was scheduled on that trial. For the “impaired outcome values” dataset,  $\alpha_{value}$  was decreased specifically for trials with inhibition scheduled during the outcome period or both periods, but not on trials with inhibition during the choice period or on control trials. For the “impaired reward processing” dataset,  $\alpha_{reward}$  was decreased on these trials instead. For the “impaired decision-making” dataset,  $\beta_{plan}$  was decreased specifically on trials following inhibition. In all cases, the parameter to be decreased was multiplied by 0.3, and synthetic datasets consisted of 100,000 total trials per rat.

### ***Histological Verification of Targeting***

We verified that surgical implants were successfully placed in the OFC using standard histological techniques. At experimental endpoint, rats with electrode arrays were anesthetized, and microlesions were made at the site of each electrode tip by passing current through the electrodes. Rats were then perfused transcardially with saline followed by formalin. Brains were sliced using a vibratome and imaged using an epifluorescent microscope. Recording sites were identified using these microlesions and the scars created by the electrodes in passing, as well dimples in the surface of the brain. Locations of optical fibers were identified using the scars created by their passage. Sharp optical fibers do not leave a visible scar near their tips, so we estimated the position of the tips using the trajectory of the scar and the known distance below brain surface to which fibers were lowered during surgery. Location of virus expression was identified by imaging the GFP conjugated to the eNpHR3.0 molecule. Note that at the time of the first submission of this paper, histological verification of targeting is still underway.

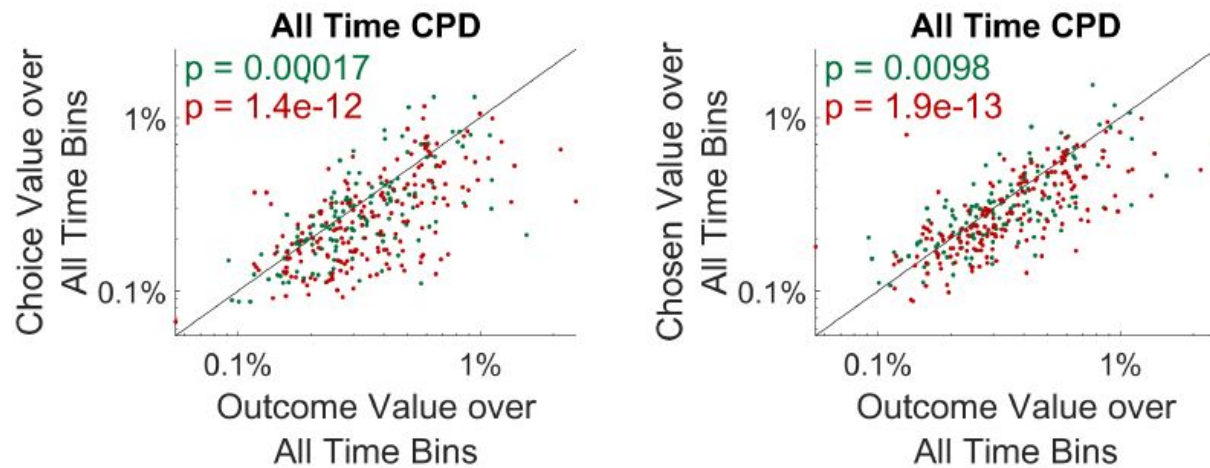
### ***Data and Code Availability***

All data collected for the purpose of this paper, and all software used in the analysis of this data, are available from the corresponding author upon reasonable request. Software used for training rats is available on the Brody lab website.

## Citations

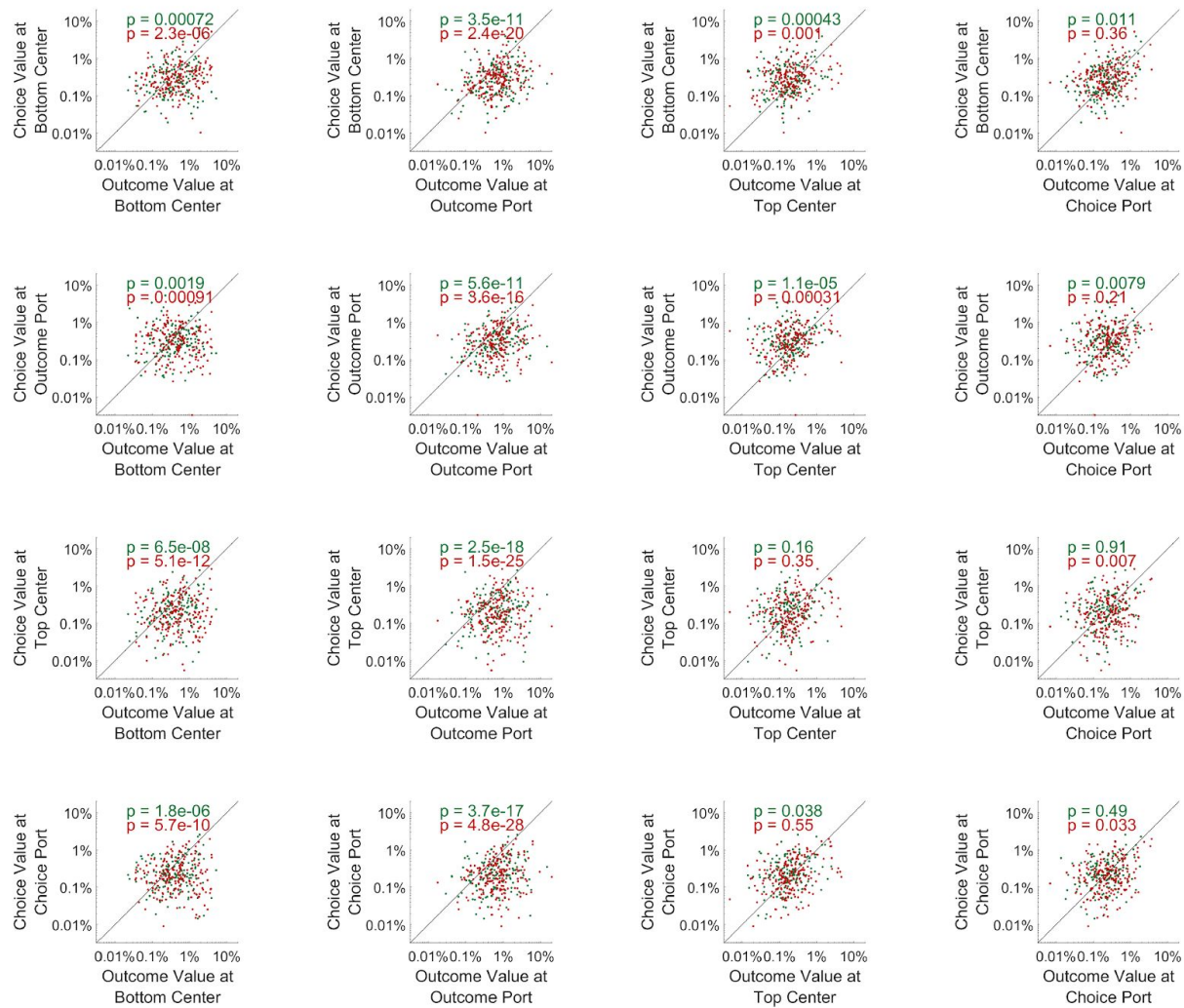
1. L. P. Sugrue, G. S. Corrado, W. T. Newsome, Choosing the greater of two goods: neural currencies for valuation and decision making. *Nat. Rev. Neurosci.* **6**, 363–375 (2005).
2. R. S. Sutton, A. G. Barto, *Reinforcement learning: An introduction* (MIT press Cambridge, 1998), vol. 1.
3. N. D. Daw, J. P. O’Doherty, in *Neuroeconomics (Second Edition)*, P. W. Glimcher, E. Fehr, Eds. (Academic Press, San Diego, 2014), pp. 393–410.
4. D. Lee, H. Seo, M. W. Jung, Neural basis of reinforcement learning and decision making. *Annu. Rev. Neurosci.* **35**, 287–308 (2012).
5. J. P. O’Doherty, Lights, camembert, action! The role of human orbitofrontal cortex in encoding stimuli, rewards, and choices. *Ann. N. Y. Acad. Sci.* **1121**, 254–272 (2007).
6. T. A. Stalnaker, N. K. Cooch, G. Schoenbaum, What the orbitofrontal cortex does not do. *Nat. Neurosci.* **18**, 620–627 (2015).
7. P. H. Rudebeck, E. A. Murray, The orbitofrontal oracle: cortical mechanisms for the prediction and evaluation of specific behavioral outcomes. *Neuron.* **84**, 1143–1156 (2014).
8. J. D. Wallis, Orbitofrontal cortex and its contribution to decision-making. *Annu. Rev. Neurosci.* (2007) (available at <http://www.annualreviews.org/doi/abs/10.1146/annurev.neuro.30.051606.094334>).
9. C. Padoa-Schioppa, Neurobiology of economic choice: a good-based model. *Annu. Rev. Neurosci.* **34**, 333–359 (2011).
10. C. Padoa-Schioppa, K. E. Conen, Orbitofrontal Cortex: A Neural Circuit for Economic Decisions. *Neuron.* **96**, 736–754 (2017).
11. L. T. Hunt *et al.*, Mechanisms underlying cortical activity during value-guided choice. *Nat. Neurosci.* **15**, 470–6, S1–3 (2012).
12. R. C. Wilson, Y. K. Takahashi, G. Schoenbaum, Y. Niv, Orbitofrontal cortex as a cognitive map of task space. *Neuron.* **81**, 267–279 (2014).
13. N. Schuck, R. Wilson, Y. Niv, A state representation for reinforcement learning and decision-making in the orbitofrontal cortex. *bioRxiv* (2017), doi:10.1101/210591.
14. A. M. Wikenheiser, G. Schoenbaum, Over the river, through the woods: cognitive maps in the hippocampus and orbitofrontal cortex. *Nat. Rev. Neurosci.* **17**, 513–523 (2016).
15. S.-L. Lim, J. P. O’Doherty, A. Rangel, The decision value computations in the vmPFC and striatum use a relative value code that is guided by visual attention. *J. Neurosci.* **31**, 13214–13223 (2011).
16. E. L. Rich, J. D. Wallis, Decoding subjective decisions from orbitofrontal cortex. *Nat. Neurosci.* **19**, 973–980 (2016).
17. G. Schoenbaum, M. R. Roesch, T. A. Stalnaker, Y. K. Takahashi, A new perspective on the role of the orbitofrontal cortex in adaptive behaviour. *Nat. Rev. Neurosci.* **10**, 885–892 (2009).
18. M. E. Walton, T. E. J. Behrens, M. P. Noonan, M. F. S. Rushworth, Giving credit where credit is due: orbitofrontal cortex and valuation in an uncertain world. *Ann. N. Y. Acad. Sci.* **1239**, 14–24 (2011).
19. H. F. Song, G. R. Yang, X.-J. Wang, Reward-based training of recurrent neural networks for cognitive and value-based tasks. *Elife.* **6** (2017), doi:10.7554/eLife.21492.
20. K. J. Miller, M. M. Botvinick, C. D. Brody, Dorsal hippocampus contributes to model-based planning. *Nat. Neurosci.* **20**, 1269–1276 (2017).
21. N. D. Daw, S. J. Gershman, B. Seymour, P. Dayan, R. J. Dolan, Model-based influences on humans’ choices and striatal prediction errors. *Neuron.* **69**, 1204–1215 (2011).
22. N. D. Daw, Y. Niv, P. Dayan, Uncertainty-based competition between prefrontal and dorsolateral striatal systems for behavioral control. *Nat. Neurosci.* **8**, 1704–1711 (2005).

23. R. J. Dolan, P. Dayan, Goals and habits in the brain. *Neuron*. **80**, 312–325 (2013).
24. D. J. Barraclough, M. L. Conroy, D. Lee, Prefrontal cortex and decision making in a mixed-strategy game. *Nat. Neurosci.* **7**, 404–410 (2004).
25. L. P. Sugrue, G. S. Corrado, W. T. Newsome, Matching behavior and the representation of value in the parietal cortex. *Science*. **304**, 1782–1787 (2004).
26. G. Corrado, K. Doya, Understanding neural coding through the model-based analysis of decision making. *J. Neurosci.* **27**, 8178–8180 (2007).
27. J. H. Sul, H. Kim, N. Huh, D. Lee, M. W. Jung, Distinct roles of rodent orbitofrontal and medial prefrontal cortex in decision making. *Neuron*. **66**, 449–460 (2010).
28. N. D. Daw, in *Decision Making, Affect, and Learning* (2011), pp. 3–38.
29. X. Cai, S. Kim, D. Lee, Heterogeneous coding of temporally discounted values in the dorsal and ventral striatum during intertemporal choice. *Neuron*. **69**, 170–182 (2011).
30. S. W. Kennerley, T. E. J. Behrens, J. D. Wallis, Double dissociation of value computations in orbitofrontal and anterior cingulate neurons. *Nat. Neurosci.* **14**, 1581–1589 (2011).
31. T. A. Stalnaker, T.-L. Liu, Y. K. Takahashi, G. Schoenbaum, Orbitofrontal neurons signal reward predictions, not reward prediction errors. *Neurobiol. Learn. Mem.* (2018), doi:10.1016/j.nlm.2018.01.013.
32. T. A. Stalnaker *et al.*, Orbitofrontal neurons infer the value and identity of predicted outcomes. *Nat. Commun.* **5**, 3926 (2014).
33. Y. K. Takahashi *et al.*, The orbitofrontal cortex and ventral tegmental area are necessary for learning from unexpected outcomes. *Neuron*. **62**, 269–280 (2009).
34. Y. K. Takahashi, T. A. Stalnaker, M. R. Roesch, G. Schoenbaum, Effects of inference on dopaminergic prediction errors depend on orbitofrontal processing. *Behav. Neurosci.* **131**, 127–134 (2017).
35. M. A. McDannald, F. Lucantonio, K. A. Burke, Y. Niv, G. Schoenbaum, Ventral striatum and orbitofrontal cortex are both required for model-based, but not model-free, reinforcement learning. *J. Neurosci.* **31**, 2700–2705 (2011).
36. M. P. H. Gardner, J. S. Conroy, M. H. Shaham, C. V. Styer, G. Schoenbaum, Lateral Orbitofrontal Inactivation Dissociates Devaluation-Sensitive Behavior and Economic Choice. *Neuron*. **96**, 1192–1203.e4 (2017).
37. G. Schoenbaum, Y. Takahashi, T.-L. Liu, M. A. McDannald, Does the orbitofrontal cortex signal value? *Ann. N. Y. Acad. Sci.* **1239**, 87–99 (2011).
38. S. C. Y. Chan, Y. Niv, K. A. Norman, A Probability Distribution over Latent Causes, in the Orbitofrontal Cortex. *J. Neurosci.* **36**, 7817–7828 (2016).
39. B. Lau, P. W. Glimcher, Dynamic response-by-response models of matching behavior in rhesus monkeys. *J. Exp. Anal. Behav.* **84**, 555–579 (2005).
40. Stan Development Team, MatlabStan: The MATLAB interface to Stan (2016), (available at <http://mc-stan.org/matlab-stan.html>).
41. Bob Carpenter, Andrew Gelman, Matt Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Michael A. Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell., Stan: A Probabilistic Programming Language (2016), (available at <http://mc-stan.org>).
42. A. Gelman *et al.*, *Bayesian Data Analysis, Third Edition* (CRC Press, 2013).
43. A. Akrami, C. D. Kopec, M. E. Diamond, C. D. Brody, Posterior parietal cortex represents sensory stimulus history and is necessary for its effects on behavior. *bioRxiv* (2017), p. 182246.
44. T. D. Hanks *et al.*, Distinct relationships of parietal and prefrontal cortices to evidence accumulation. *Nature*. **520**, 220–223 (2015).
45. C. D. Kopec, J. C. Erlich, B. W. Brunton, K. Deisseroth, C. D. Brody, Cortical and Subcortical Contributions to Short-Term Memory for Orienting Movements. *Neuron*. **88**, 367–377 (2015).

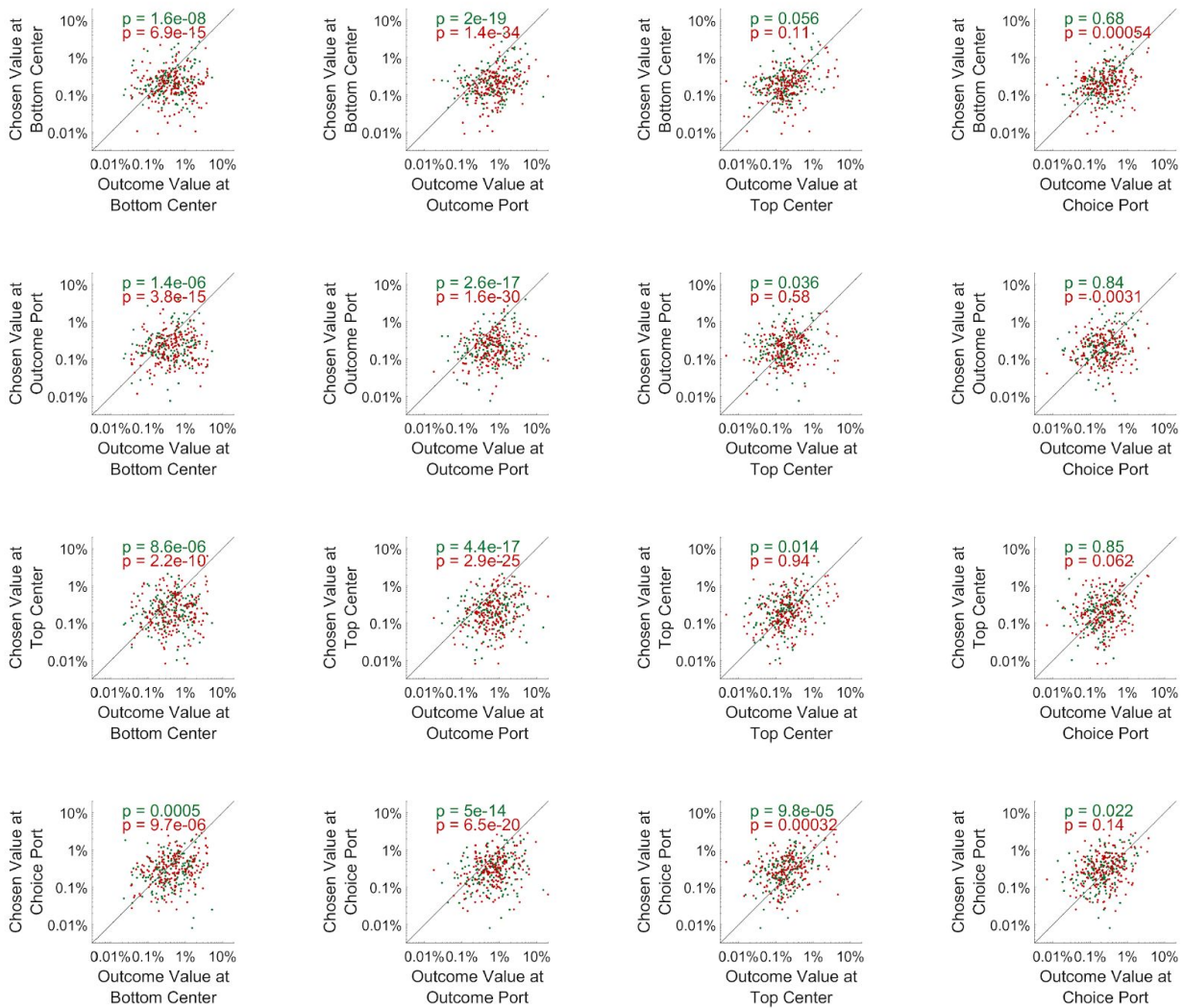


**Figure ED1:** Left: Coefficient of Partial Determination (CPD) for the outcome-value and the choice-value-difference regressors, computed aggregating variance over all time bins for each single unit (green) and multi-unit (red) cluster. P-values shown are the result of a sign test over units. Right: CPD for the outcome-value and the chosen-value predictors.

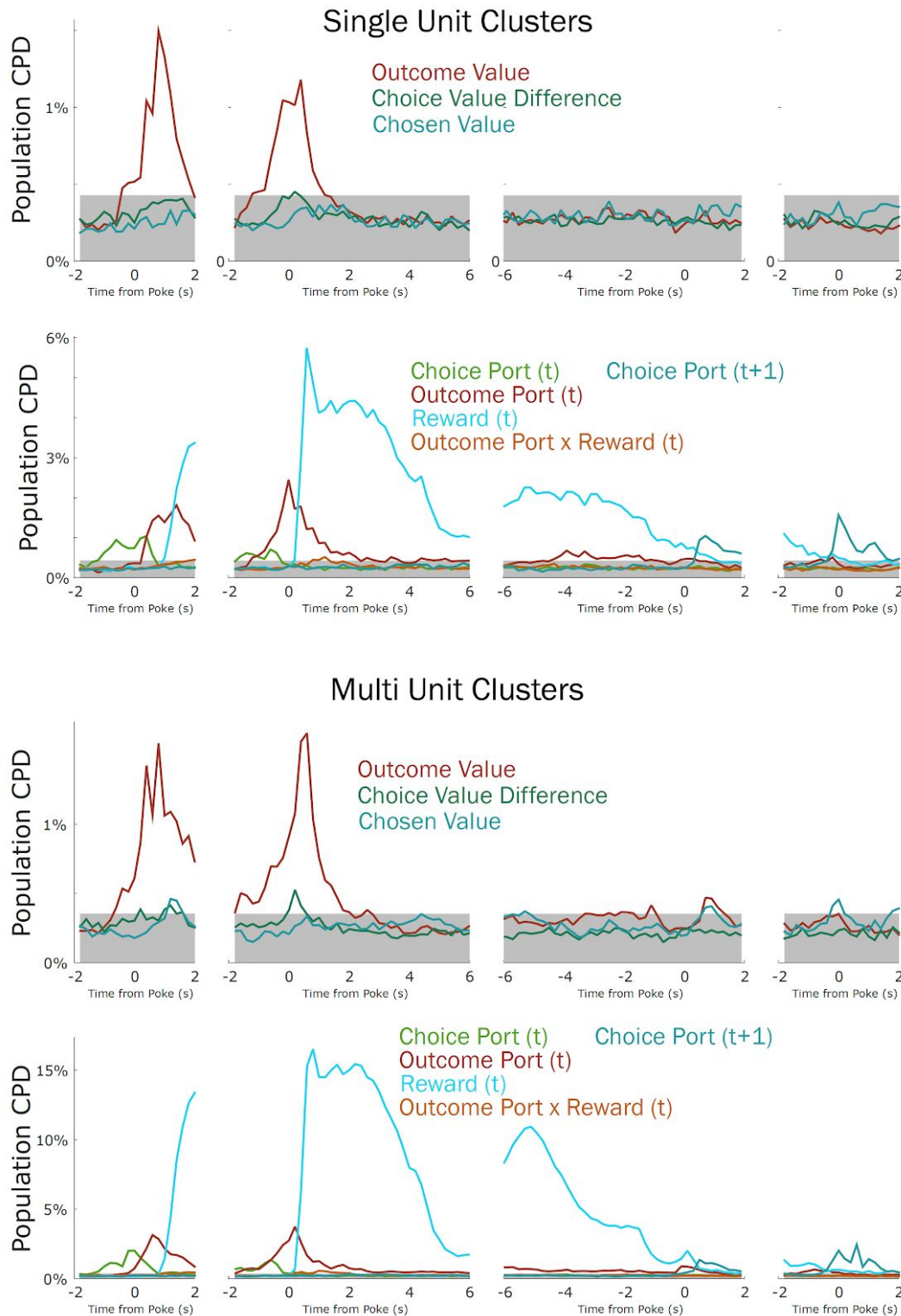




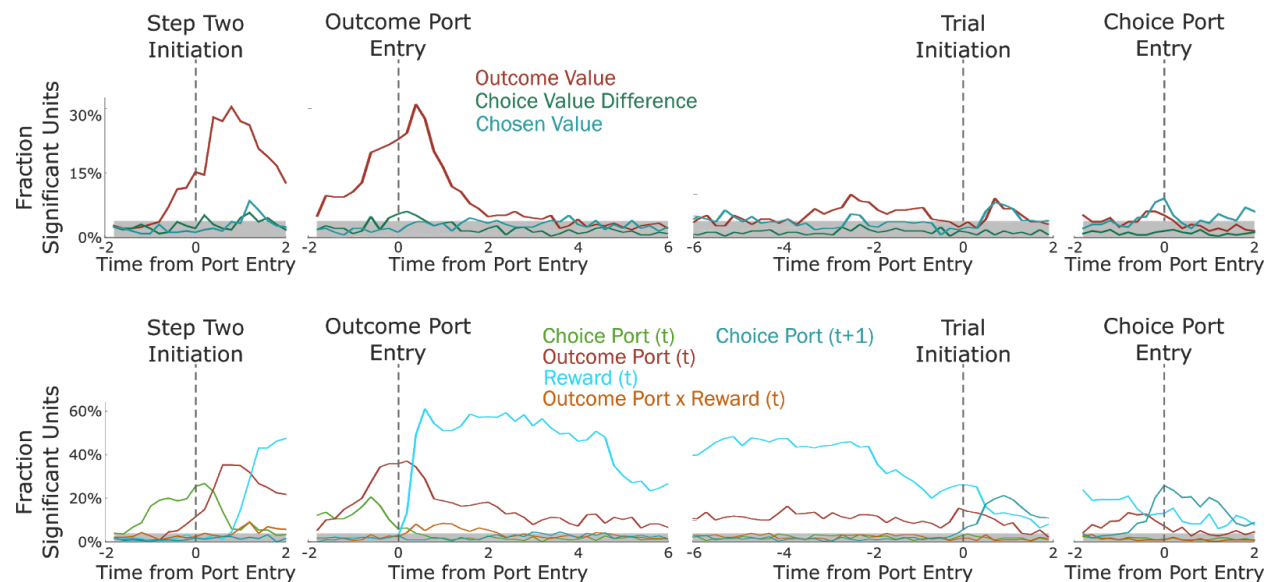
**Figure ED2:** Coding of outcome value and choice value difference at the time of port entry. Each panel shows CPD for the outcome-value and the choice-value-difference regressor, each computed in a one-second window (five 200ms time bins) centered on a different port entry event. P-values shown are the result of a sign test over units.



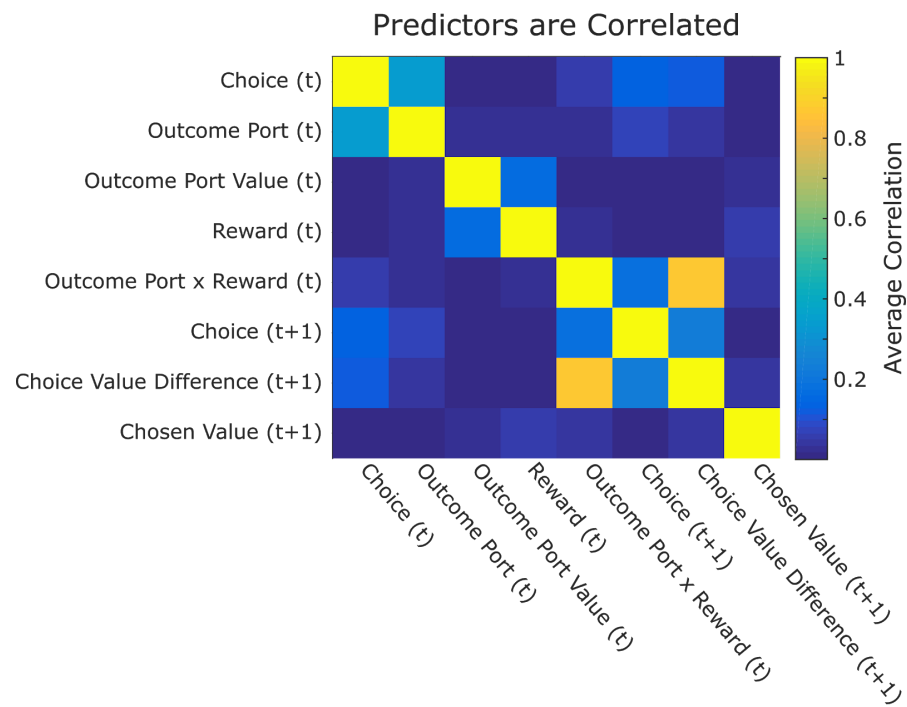
**Figure ED3:** Coding of outcome value and chosen value at the time of port entry. Each panel shows CPD for the outcome-value and the choice-value-difference regressor, each computed in a one-second window (five 200ms time bins) centered on a different port entry event for single-unit (green) and multi-unit (red) clusters. P-values shown are the result of a sign test over units.



**Figure ED4:** Timecourse of population CPD for the six predictors in our model (see also **Fig. 2a**), considering only single-unit clusters (above) or considering only multi-unit clusters (below)

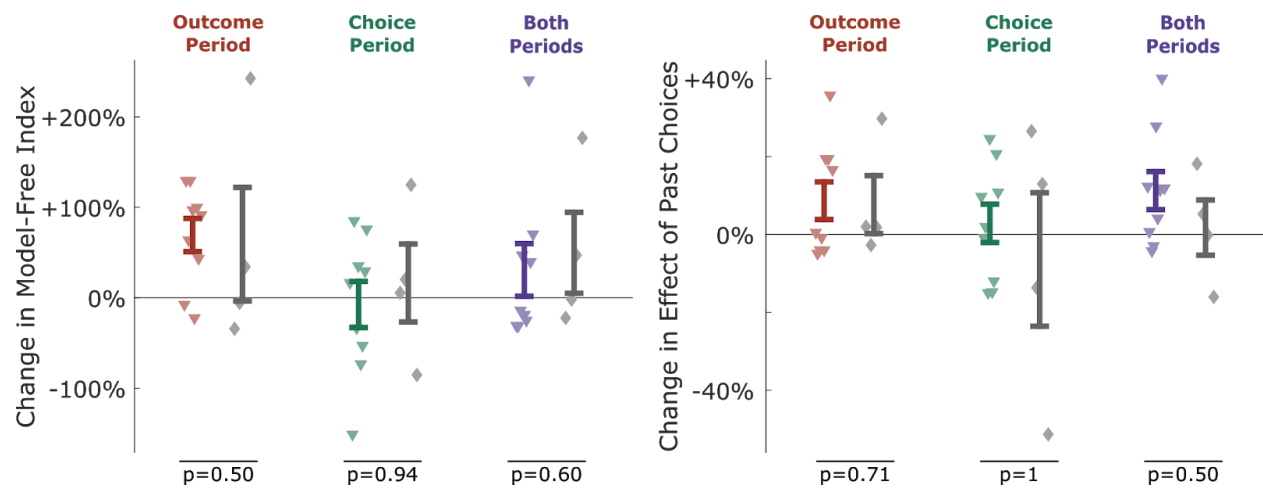


**Figure ED5:** Fraction of units significantly encoding each predictor in each 200ms time bin. Units were deemed significant in a bin for a predictor if they earned a coefficient of partial determination larger than that of 99% of permuted datasets for that predictor in that bin. This plot includes both single- and multi-unit clusters.

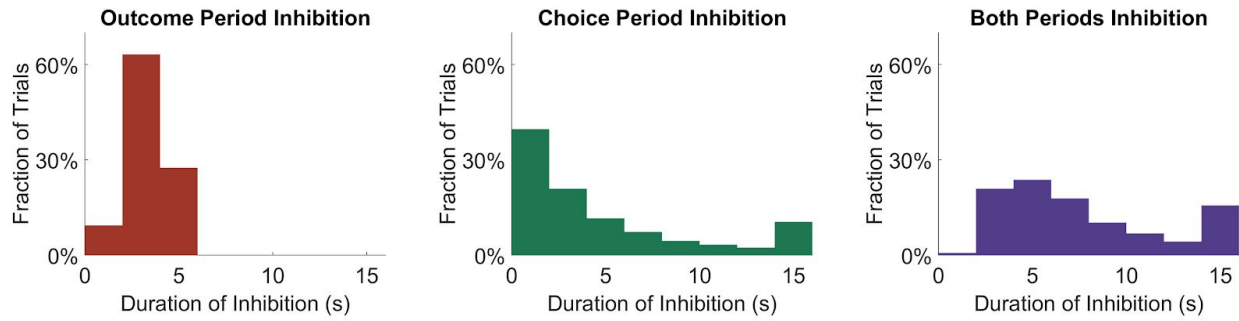


**Figure ED6:** Correlations among predictors in the model





**Figure ED7:** Effects of optogenetic inhibition on the model-free index and on the main effect of past choices. No significant differences were found between inhibition and sham inhibition on either of these measures (rank sum tests, all  $p > 0.5$ ).



**Figure ED8:** Duration of inhibition associated with outcome-period, choice-period, or both-periods conditions.