Behavioral/Cognitive

# A Probability Distribution over Latent Causes, in the Orbitofrontal Cortex

Stephanie C. Y. Chan, Yael Niv,* and Kenneth A. Norman*

Princeton Neuroscience Institute, Princeton University, Princeton, New Jersey 08544

The orbitofrontal cortex (OFC) has been implicated in both the representation of "state," in studies of reinforcement learning and decision making, and also in the representation of "schemas," in studies of episodic memory. Both of these cognitive constructs require a similar inference about the underlying situation or "latent cause" that generates our observations at any given time. The statistically optimal solution to this inference problem is to use Bayes' rule to compute a posterior probability distribution over latent causes. To test whether such a posterior probability distribution is represented in the OFC, we tasked human participants with inferring a probability distribution over four possible latent causes, based on their observations. Using fMRI pattern similarity analyses, we found that BOLD activity in the OFC is best explained as representing the (log-transformed) posterior distribution over latent causes. Furthermore, this pattern explained OFC activity better than other task-relevant alternatives, such as the most probable latent cause, the most recent observation, or the uncertainty over latent causes.

*Key words:* Bayes' rule; context; posterior distribution; schemas; state representation; ventromedial prefrontal cortex

---

**Significance Statement**

Our world is governed by hidden (latent) causes that we cannot observe, but which generate the observations we see. A range of high-level cognitive processes require inference of a probability distribution (or "belief distribution") over the possible latent causes that might be generating our current observations. This is true for reinforcement learning and decision making (where the latent cause comprises the true "state" of the task), and for episodic memory (where memories are believed to be organized by the inferred situation or "schema"). Using fMRI, we show that this belief distribution over latent causes is encoded in patterns of brain activity in the orbitofrontal cortex, an area that has been separately implicated in the representations of both states and schemas.

---

## Introduction

In recent years, cognitive neuroscientists studying reinforcement learning and decision making have recognized the importance of specifying representations of the environmental "state" that capture the structure of the world in a predictive way (Courville et al., 2004; Gershman and Niv, 2010). At the same time, there has been renewed interest among cognitive neuroscientists in how memory encoding and retrieval are shaped by situation-specific prior knowledge ("schemas"; Tse et al., 2007). As work in this area

progresses, it is important to clarify exactly what constitutes a schema and how schemas are formed.

Whether inferring the current state or the currently relevant schema, agents are making inferences about the hidden variables that underlie and generate our observations in the world. This inference can be concretely formulated in terms of Bayesian latent cause models (e.g., Gershman et al., 2010). According to this framework, states and schemas can be viewed as hidden (latent) causes that give rise to observable events. For example, if you arrive late to a lecture, the situation (whether this is indeed the department colloquium or you have accidentally walked in on an undergraduate class) determines your observations about the average age of the audience, the proportion of audience members that are taking notes, the type of information being presented, and so on. To decide whether you are in the right place, you can use Bayesian inference to infer a belief distribution over the possible situations that might have generated the current observations, i.e., a posterior probability distribution over latent causes, $p(latent\ cause\ |\ observations)$ (Fig. 1A).

We hypothesized, based on the similarity of the underlying computations, that the inference related to these two cognitive constructs
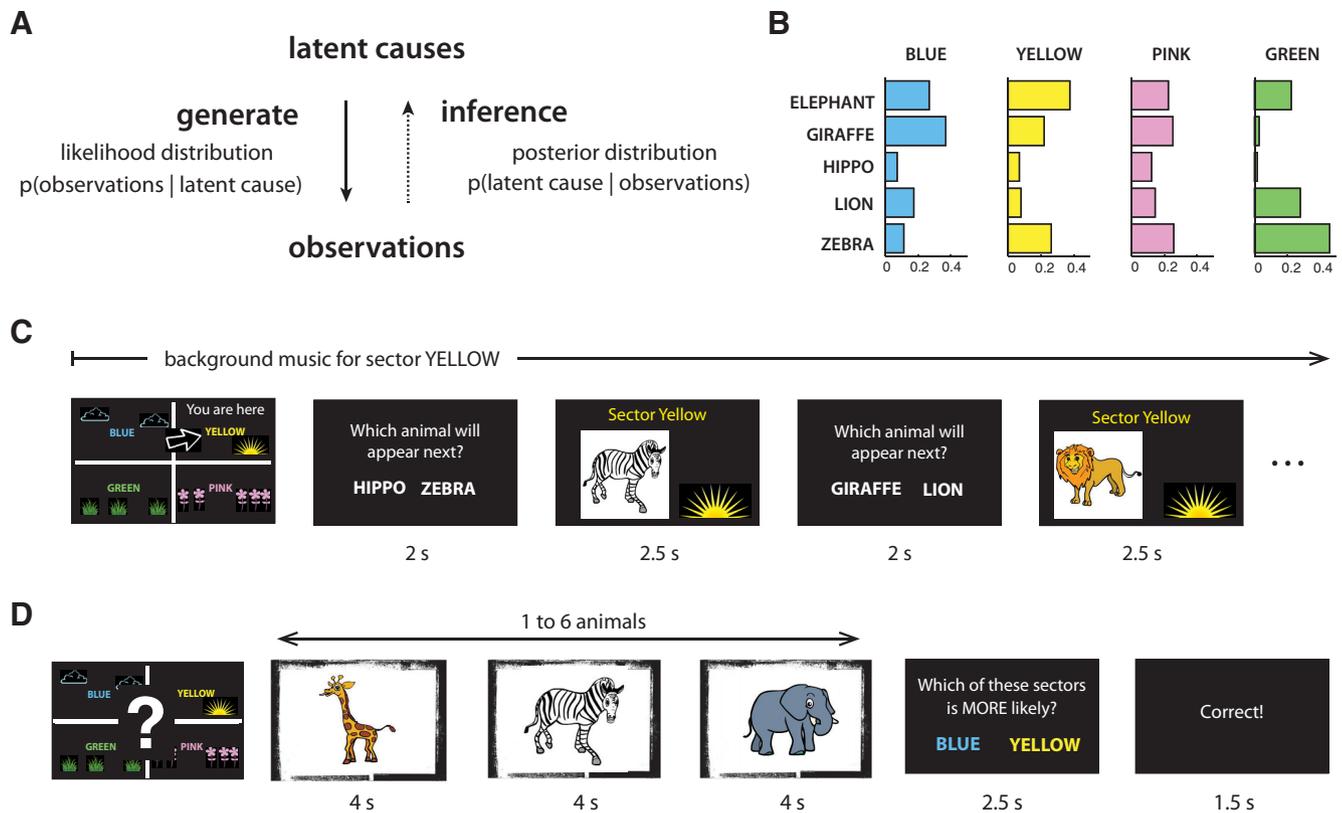
**Figure 1.** Task. ***A***, Schematic showing the relationship between latent causes and observations in the world. Inference about the posterior probability over latent causes involves inverting the generative model. ***B***, Animal likelihood distributions *P*(*animal* | *sector*) (not shown directly to participants). Colors and animals were randomized across participants. ***C***, An example of the first few trials of a tour through sector YELLOW. Each tour began with an image of the safari map, indicating the current sector and its location, and lasted 30 – 40 trials. Each trial began with a prompt asking the participant to guess which animal would appear next, followed by the appearance of an animal. A fixation cross was presented for 0.2– 0.8 s before each question and each animal presentation. The animals were pseudorandomly drawn from the likelihood distributions for the current sector. The sector's music played in the background, until the start of the next tour. ***D***, An example of a trial in the photographs task. Each trial began with an image of the safari map with a question mark at its center, indicating that the current sector was unknown. Next, a sequence of 1– 6 animals appeared (pseudorandomly drawn from a single sector). Finally, participants were prompted to guess which of two sectors (randomly chosen) was more (or, on half the trials, less) probable. Participants received feedback on their responses. A fixation cross was presented for the last 0.5 s of each animal presentation.

(states and schemas) might be implemented using the same neural hardware. Indeed, there is one area of the brain that has separately been implicated in representing states (Wilson et al., 2014) and also schemas (Tse et al., 2011; Ranganath and Ritchey, 2012; van Kesteren et al., 2012; Ghosh and Gilboa, 2014; Richards et al., 2014; Schlichting and Preston, 2015)—the orbitofrontal cortex (OFC). Furthermore, previous univariate analyses in fMRI have implicated this region in encoding various summary statistical measures that are related to or are components of the posterior distribution, e.g., the posterior mean, likelihood of the current stimulus, and prior uncertainty (Vilares et al., 2012; d'Acremont et al., 2013; Ting et al., 2015). However, these studies have not investigated representations of a full probability distribution.

Here, we used fMRI to investigate representation in the OFC of posterior probability distributions over latent causes. In our experiment, we created a probabilistic environment in which participants were required to make inferences about the hidden causes that generated their observations. Participants viewed sequences of animal "photographs" from one of four "sectors" of a virtual "animal reserve." Participants were asked to judge the probability that the photographs were taken in each of the sectors, based on their previous experience observing animals in each sector. Using multivariate pattern similarity analyses of fMRI activity, we found that BOLD activity in the OFC was better explained by the posterior distribution over sectors (latent

causes) than by a wide range of related signals, including the current stimulus, the most probable sector (the maximum a posteriori latent cause), or the uncertainty over latent causes (operationalized as the entropy of the posterior distribution). The present result advances our understanding of the function of the OFC. It also unifies results from two different fields of cognitive neuroscience, inviting further investigation into the relationship between probabilistic inference, states, and schemas.

## Materials and Methods

### Participants
Thirty-two participants (age, 18–34 years; 22 female) from the Princeton University community participated in exchange for monetary compensation ($20 per hour plus up to $15 performance-related bonus). All participants were right-handed. Participants provided informed written consent. The study was approved by the Princeton University Institutional Review Board.

### Experimental design
*The safari.* Participants were told that they were going on a virtual safari—a visit to an animal reserve divided into four different sectors. Each sector was associated with a different color, background image, background music, and location on a map (randomized across participants).

There were five different kinds of animals in the animal reserve. Every animal appeared in every sector, but with different likelihoods

**Table 1. Tasks performed by participants on Day 1 and Day 2**

| Day 1 | | |
|---|---|---|
| Tours task–1st set | 40 trials each tour | 2 tours through each sector, going clockwise around the map |
| Tours task–2nd set | 20 trials each tour | 2 tours through each sector, sectors pseudorandomly ordered |
| Day 2 | | |
| Tours task–3rd set | 30 trials each tour | 2 tours through each sector, sectors pseudorandomly ordered |
| Tours task–4th set | 10 trials each tour | 2 tours through each sector, sectors pseudorandomly ordered |
| Photographs task–1st set | 2 sessions × 20 trials each | Outside the MRI scanner |
| Photographs task–2nd set | 4 sessions × 30 trials each | Inside the MRI scanner |

P(*animal* | *sector*). The likelihoods (not shown directly to the participants) were chosen so that none of the sectors was strongly identified with a single animal, and so that none of the animals was strongly identified with a single sector (Fig. 1B; colors and animals were randomly assigned across participants).

*Procedure overview.* The experiment consisted of two parts. In the first part, participants "toured" through the animal reserve to learn (through experience) the likelihoods P(*animal* | *sector*) for each animal and each sector. In the second part of the experiment, participants were shown sequences of drawings that represented "photographs" of animals, which were taken in an unknown sector. Given the animals shown in each sequence, participants were asked to infer the posterior probabilities of different sectors, P(*sector* | *animals shown*). For each participant, the experiment took place across 2 consecutive days (Table 1).

*Tours task.* In the tours task (Fig. 1C), participants were instructed that they would "tour" through the animal reserve, one sector at a time, to learn the animal frequencies in each sector (the animal likelihoods). Each tour consisted of 10–40 trials through a single sector. One animal appeared on each trial, pseudorandomly chosen according to the likelihoods for that sector. Before each animal appeared, participants were shown a prompt, asking them to make a prediction about which of two animals (one correct and one randomly chosen) would appear next. The alternative (incorrect) option was chosen with uniform probability from the four other animals. To distinguish between the animals in the question prompt (which were not representative of the sector's likelihood distribution) and the animals that were actually drawn from the sector's likelihood distribution, the question prompts were shown as text while the animals drawn from the safari sector were shown as pictures.

In order for the sectors to form rich contexts, each sector was associated with a different color, background image, background music, and location on a two-by-two map (randomized across participants). Before the first trial of a tour through a sector, participants were shown the sector's location on the map. Also, for the duration of a tour through a sector, animals were displayed on the sector's color-matched backdrop image, and the music associated with that sector was played in the background.

*Photographs task.* On each trial of the photographs task (Fig. 1D), participants were shown a sequence of animal "photographs," without being told which sector the photographs were taken from. At the end of the sequence, participants were prompted to indicate which of two sectors (randomly chosen) was more (or less) probable. The two sector options for each question were chosen uniformly from the four sectors of the safari (and did not necessarily include the most or least likely sector). So, to perform well on the task, participants had to maintain a full posterior distribution over all four sectors (as opposed to estimating only the most probable sector, for instance).

Participants received 10 cents for every correctly answered question, and they received feedback on every trial. So that more probable sectors were not consistently associated with higher monetary value, we asked

which of the two sectors was more probable on half of the trials, and which was less probable on the other half of the trials. To eliminate confounds with motor plan, the positions of the two response options were pseudorandomly assigned between left and right.

To encourage participants to update their inference of the sector probabilities after every animal presentation (as opposed to waiting until the time of the question to integrate over the animals observed), we varied the length of the sequences between one and six animals (so that the appearance of the question prompt was unpredictable), and participants were only allowed 2.5 s to give a response after the appearance of the question.

The posterior probability of each sector P(*sector* | *animals seen*) can be straightforwardly computed from the animal likelihoods using Bayes' rule (all sectors were equally likely a priori), as follows (Eq. 1):

$$P(\text{sector } i | \text{animals seen}) \propto P(\text{animals seen} | \text{sector } i) \times P(\text{sector } i)$$

$$\propto P(\text{animal 1} | \text{sector } i) \times P(\text{animal 2} | \text{sector } i) \times \dots$$

Feedback for the responses was generated based on these posterior probabilities. Due to a bug in the code that was undetected during data collection, the feedback was incorrectly generated for some of the trials containing only one animal presentation (this affected ~10% of the trials). In our fMRI analyses, to account for learning from the incorrect feedback, we used each participant's estimates of the likelihoods (collected at the end of the experiment) instead of the real likelihoods, and we also performed trial-by-trial behavioral model fitting, in order to model learning from feedback (see next section).

To familiarize themselves with the photographs task, participants first performed two sessions (20 trials each) of the task outside the MR scanner. They then performed four sessions (30 trials, ~11 min per session) inside the scanner.

*Behavioral model fitting*

To model participants' posterior inference on the photographs task, as well as any learning from feedback, we performed trial-by-trial model fitting of participants' responses. We tested several classes of behavioral models (note that we will explicitly refer to these as "behavioral models," to distinguish them from the neural models that we later test against the neural data). These classes of models (Bayesian_nolearning, Additive, Most/least voter, and Bayesian_feedbackRL) are described in the following.

*Bayesian_nolearning.* This behavioral model assumed that participants were correctly computing the posterior distribution over sectors P(*sector* | *animals seen*) using Bayesian inference (as in Eq. 1). To obtain the model-derived likelihood of each behavioral response (and to capture stochasticity in participants' behavior), we used a softmax on the posterior probabilities of the two options in each question prompt (Eq. 2):

$$P(\text{response} = \text{option 1})$$

$$= \frac{1}{1 + \exp[\beta * (P(\text{sector} = \text{option 1}|\text{animals seen}) - P(\text{sector} = \text{option 2}|\text{animals seen}))]}$$

where $\beta$ is an inverse temperature parameter ($\beta = 0$ implies equal likelihood for both options).

*Additive.* In this behavioral model, rather than correctly multiplying the animal likelihoods together to obtain the posterior distribution over sectors (as in Eq. 1), we assumed that participants instead added the likelihoods together to obtain an "additive posterior" (normalized to sum to 1; Eq. 3):

"Additive posterior" $\propto P(\text{animal 1}|\text{sector}) + P(\text{animal 2}|\text{sector}) + \dots$

While statistically suboptimal, we might expect this from a simple associative mechanism that brings the sectors to mind in proportion to their association strength with the animals seen. Again, to determine response probabilities, we applied a softmax operator to the additive "posterior" probabilities for the two options in each question prompt.

*Most/least voter.* These behavioral models assumed that participants were only paying attention to the most common (and/or least common)

animals in each sector, a similar strategy having been previously observed in a similar task (Gluck et al., 2002). During the trials, each animal appearance "voted" for (or against) the sectors in which it was the most common (or least common). To obtain the model-derived likelihood of each behavioral response, we used a softmax on the final tally at the end of each sequence.

We tested several variants of this behavioral model, e.g., tallying only the positive votes, and/or allowing an animal to "vote" for (or against) a sector if it was one of the *two* most (or least) common animals in that sector. The magnitude of the positive and negative votes were either allowed to be two separate free parameters or constrained to be equal to each other. Because the magnitude of the vote already served as a scaling parameter for the input to the softmax operator, the inverse temperature of the softmax was kept constant at 1.

*Bayesian_feedbackRL.* These behavioral models were designed to account for learning from feedback during the photographs task (including the incorrectly generated feedback). Here we assumed a reinforcement learning process, in which participants adjusted their internal estimates of the animal likelihoods after feedback about the two sectors in the question. These likelihoods were then used to compute the posterior distribution over sectors via Bayes' rule.

For the sector that feedback indicated to be more probable, likelihoods were adjusted upward for all animals seen on that trial. For the sector indicated to be less probable, likelihoods were adjusted downward for all animals seen on the trial (Fig. 2; Eq. 4):

$$P(\text{animal}|\text{more probable sector})_{new}$$

$$= P(\text{animal}|\text{more probable sector})_{old}$$

$$+ \alpha_{pos}(1 - P(\text{animal}|\text{more probable sector})_{old})$$

$$P(\text{animal}|\text{less probable sector})_{new}$$

$$= (1 - \alpha_{neg})P(\text{animal}|\text{less probable sector})_{old}$$

Estimates of the likelihoods were renormalized after each adjustment. The learning rates $\alpha_{pos}$ and $\alpha_{neg}$ were either allowed to be two separate free parameters or they were constrained to be equal.

For the initialization of the likelihoods, we tested two versions: initialization at the true animal likelihoods or initialization according to the participants' subjective estimates of the likelihoods (collected at the end of the experiment; see below).

Finally, the likelihoods were used to compute the posterior distribution over sectors via Bayes' rule. Thus, posterior inference in the FeedbackRL behavioral model also used Bayes' rule. The only difference from the "Bayesian_nolearning" behavioral model above is that the likelihoods (which enter into the posterior inference computation from Eq. 1) were adjusted on each trial according to feedback.

We tested several additional variants of this behavioral model. In one variant, participants only adjusted their likelihoods in response to "You are incorrect" feedback (instead of in response to all feedback). In another variant of the behavioral model, we scaled the learning rates separately for each animal according to how much that animal contributed to the final posterior distribution (Eq. 5):

$$\alpha_{eff,\text{animal X}}$$

$$= \alpha \cdot abs[P(\text{more probable sector}|\text{appearances of animal X})$$

$$- P(\text{less probable sector}|\text{appearances of animal X})]$$
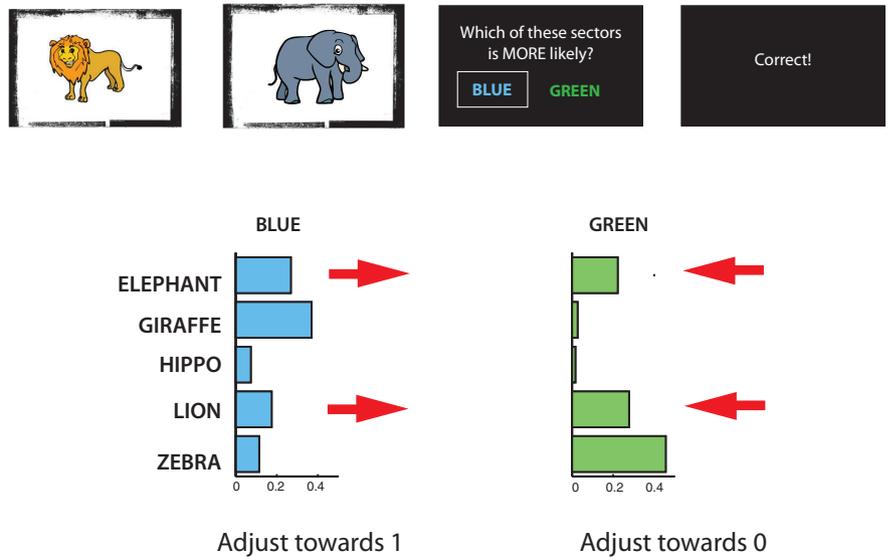


**Figure 2.** FeedbackRL behavioral model. An illustration of learning from feedback in the Bayesian_feedbackRL behavioral model, for a single trial (not real data). In this example trial, the participant saw a lion and an elephant, and was asked about sector BLUE and sector GREEN. The feedback indicated that sector BLUE was more probable. As a result, the likelihoods $P(BLUE | elephant)$ and $P(BLUE | lion)$ are adjusted toward 1 with learning rate $\alpha_{pos}$, and the likelihoods $P(GREEN | elephant)$ and $P(GREEN | lion)$ are adjusted toward 0 with learning rate $\alpha_{neg}$.

In this variant, animals appearing multiple times in a trial would have higher effective learning rates, having contributed more to the final decision.

*For all behavioral models.* In a postexperiment questionnaire, we asked participants to provide their estimates for the animal likelihoods in each sector. For each of the behavioral models above, we tested versions using (1) the actual animal likelihoods and (2) subjective estimates of the animal likelihoods. For the few participants who provided likelihood estimates that did not sum to 1, we normalized the estimates. To avoid taking logarithms of 0, we converted estimated likelihoods of 0 into 0.01 (and renormalized).

For each of the behavioral models, we also tested versions in which the earlier and/or later animals in each sequence were given extra weight. To model these primacy/recency effects, we fit a power law function for each participant to give more weight to the earlier and/or later animals in each sequence (e.g., $1^w, 2^w, \ldots$ for animal 1, animal 2, ...). The likelihoods were exponentiated by this weighting and renormalized. If modeling both recency and primacy, the weightings for each were summed. We tested versions in which the recency and primacy free parameters $w$ were either allowed to be two free parameters or they were constrained to be equal.

Free parameters for each behavioral model were fit to each participant's behavioral data separately, using Matlab's "fmincon" function, with ≥10 random initializations for each behavioral model and each participant. The best-fitting parameters (the maximum likelihood estimates) were used to evaluate, for each participant and each behavioral model, the (geometric) mean likelihood per trial (i.e., the exponentiated log likelihood per trial, without any penalization for number of parameters), the Akaike information criterion (AIC), and the Bayesian information criterion (BIC), to compare the behavioral models and determine which best accounted for participants' behavior.

*fMRI acquisition and preprocessing*
Functional brain images were acquired using a 3 T MRI scanner (Skyra, Siemens) and preprocessed using FSL [Functional MRI of the Brain (FMRIB) Software Library; http://fsl.fmrib.ox.ac.uk/fsl/]. An echoplanar imaging sequence was used to acquire 36 slices (3 mm thickness with 1 mm gap; TR = 2 s; TE = 27 ms; flip angle, 71°). To increase signal in the OFC, slices were angled ~30° from the axial plane toward a coronal orientation (Deichmann et al., 2003). For each participant, there were four scanning runs in total (~11 min each). The functional images were
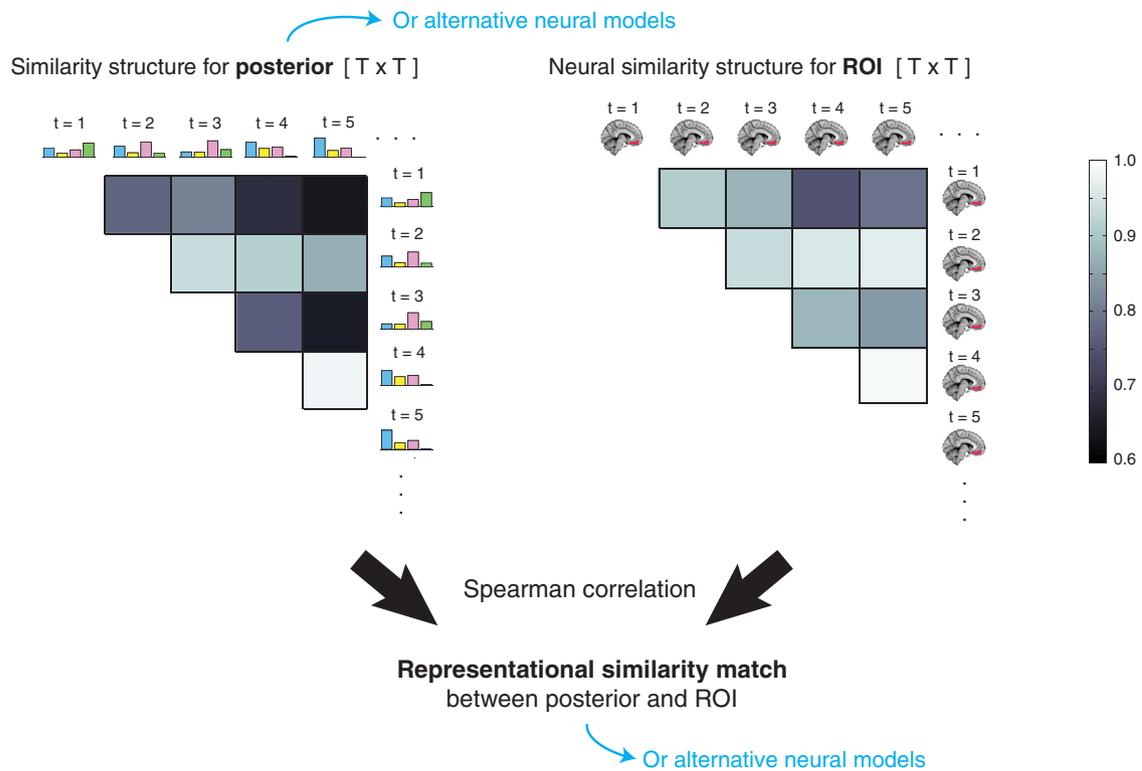
**Figure 3.** Representational similarity analysis. An illustration of the representational similarity analysis (not real data). We first computed the similarity structure for the posterior distribution (or any alternative neural model; Table 3) by computing the normalized correlation of the posterior at every time point with every other time point. We also computed the neural similarity structure for our ROI (or for each searchlight in the whole-brain analysis), by computing the normalized correlation between patterns of activity at every time point with every other time point. To evaluate the representational similarity match between the neural data and the neural model, we then computed the Spearman correlation between the two matrices (using only the upper triangle of each matrix, excluding the diagonal).

spatially filtered using a Gaussian kernel (full width at half maximum of 5 mm), and temporally filtered using a high-pass cutoff of 0.0077 Hz. We performed motion correction using a six-parameter rigid body transformation to coregister functional scans, and then registered the functional scans to an anatomical scan using a six-parameter affine transformation.

The motion regressors (and their derivatives) were residualized out from the functional images, as were the mean time courses for CSF and white matter [segmentation was performed using FSL's FAST (FSL Automated Segmentation Tool) function], and also the mean time course for blood vessels (estimated by taking voxels with the top 1% in SD across time). Then, the functional images were $z$-scored over time. All analyses were performed for each participant in participant space, and then spatially normalized by warping each participant's anatomical image to MNI space using a 12-parameter affine transformation.

*Region of interest—suborbital sulcus*
Our region of interest (ROI) was determined as the intersection of two sets of brain areas. The first set of areas, the OFC, has been postulated to be involved in the representation of state, due to evidence from studies of human and animal reinforcement learning and decision making (Wilson et al., 2014). The second set of areas, sometimes referred to as the "posterior medial network," has been postulated to be involved in the computation and representation of schemas or context (Ranganath and Ritchey, 2012), as the set of areas with high connectivity with parahippocampal cortex (PHC). The intersection of these sets of areas is the suborbital sulcus, a medial subregion of the OFC (see Fig. 6A). Using Freesurfer (Destrieux et al., 2010), the ROI was drawn as the anatomically parcellated cortical region centered on the voxel with maximal resting-state functional connectivity to the PHC (Libby et al., 2012; the ROI comprised 105 voxels in MNI 3 mm space and 97.3 ± 2.6 voxels in subject space).

*Representational similarity analysis*
If the suborbital ROI contains a multivariate representation of the posterior distribution over latent causes, then patterns of neural activity in this area should be more similar for pairs of time points at which the posterior distribution was similar, and they should be dissimilar for pairs of time points at which the posterior distribution was dissimilar. Therefore, to test whether multivariate patterns of activity in the ROI might be representing the posterior distribution over sectors, we performed a representational similarity analysis (Kriegeskorte et al., 2008).

We first computed the similarity of the posterior distribution over sectors for every pair of time points during which we expected the posterior distribution to be updated (i.e., at the times of the animal appearances). This provided us with the similarity matrix for the posterior. We also computed the similarity of the neural pattern in the ROI for every pair of time points—the similarity matrix for the ROI. Then we computed the Spearman rank correlation of these two matrices (taking only the upper triangle and excluding the diagonal). We denote this Spearman correlation as the similarity match between the posterior and the ROI (Fig. 3). We expected the similarity match to be positive, i.e., that the neural patterns in the ROI should be more similar for pairs of time points at which the posterior distribution over sectors was more similar.

We also computed the similarity match for the ROI with other signals (henceforth called "neural models"), to compare with the similarity match between the ROI and the posterior distribution over latent causes. This is important because the similarity structure for the ROI could potentially be correlated with the similarity structure of the posterior distributions for reasons other than that the posterior distribution is represented in this area. For example, the posterior distribution is, on average, more similar for pairs of time points at which the same animal is presented; if the suborbital ROI represents the animal currently presented, we would also find a positive similarity match between the ROI and the poste-
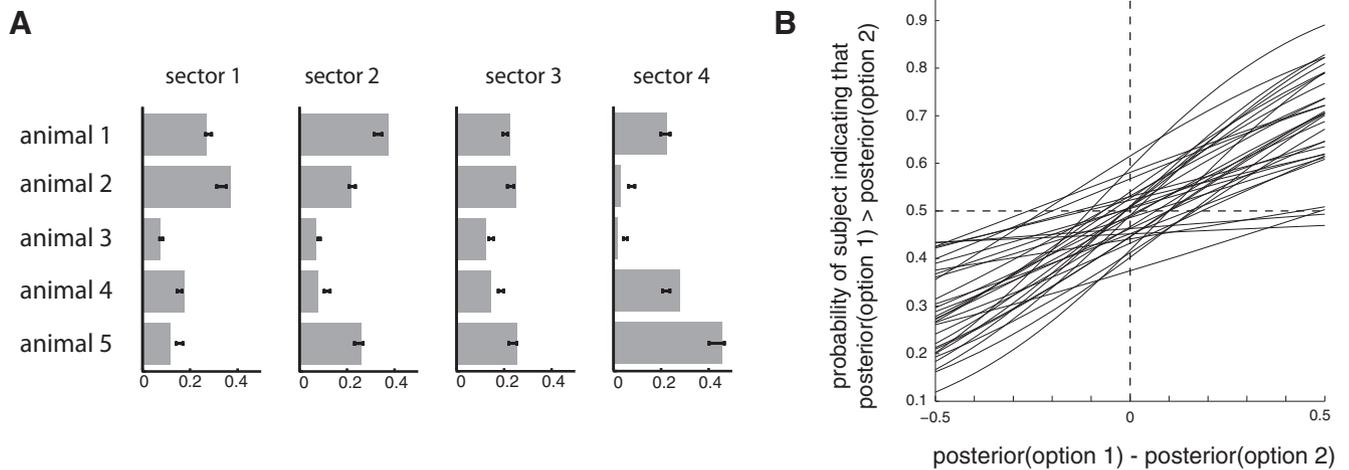
**Figure 4.** Behavioral performance. ***A***, Participants' subjective estimates of the animal likelihoods *P(animal | sector)*, for each animal and each sector, collected in a postexperiment questionnaire. Gray bars indicate the true likelihoods, black intervals indicate the mean estimates ± SEM. ***B***, Logistic regression on participants' responses during the fMRI scan sessions suggests that participants learned and used the full posterior distributions (each line shows logistic regression for one participant). The *x*-axis indicates the difference in posterior probability between the first and second options in the question. The *y*-axis indicates the probability that a participant would indicate that the first option has higher posterior probability than the second option. Mean regression parameters across participants: slope, 1.8 ± 0.040; intercept, −0.04 ± 0.15.

rior distribution. We therefore compared the similarity match between the ROI and each alternative neural model, to determine the neural model that best explained the similarity structure of the neural data.

The set of alternative neural models used for this comparison included the log-transformed posterior distribution (because many signals in the brain are known to be represented in log space; e.g., Gibbon, 1977; Yang and Shadlen, 2010; Longo and Lourenco, 2007), the current stimulus, the maximum a posteriori (MAP) sector (most probable sector), the entropy of the posterior distribution (a proxy for overall uncertainty), the probability of the MAP sector (a proxy for confidence, acting approximately as the converse of the entropy), the rank ordering of sectors in the posterior (since the task tests participants' ability to rank order the sectors), and temporal distance between measurements (because fMRI pattern similarity is known to vary as a function of the temporal distance between measurements). We also included neural models of the posterior and MAP that were instead derived using the Bayesian_feedbackRL behavioral model (given that this was the best behavioral model after Bayesian_nolearning, as determined from behavioral model fitting, described above). See Table 3 for a full list of neural models tested.

To investigate the specificity of the result to our ROI, we also performed a whole-brain "searchlight" analysis, using 25-voxel spherical searchlights. As with the ROI, we computed the similarity of the neural patterns in each searchlight to obtain the neural similarity matrix for the searchlight. We then computed the Spearman correlation of the similarity matrix for each searchlight with each of our neural models. The analysis was repeated for a searchlight centered on every voxel in the brain.

For both the ROI and searchlight analyses, the neural pattern for each animal appearance was averaged over the two TRs during which the animal appeared on the screen (after correcting for the hemodynamic lag with a 4 s shift). Similarity for neural patterns was computed using normalized correlation, to accord with the similarity measure used for the posterior-based neural models (similar results are obtained when using Pearson correlation instead). Searchlight results are displayed on an inflated brain, using the AFNI (Analysis of Functional NeuroImages) SUMA (Surface Mapping with AFNI) surface mapper (http://afni.nimh.nih.gov/afni/suma).

*Statistics and confidence intervals*
Unless stated otherwise, all statistics were computed using random-effects bootstrap distributions on the mean by resampling participants with replacement (Efron and Tibshirani, 1986). All confidence intervals in the text are given as SEM.

To test the reliability of searchlight results across participants, we used the "randomize" function in FSL (http://fsl.fmrib.ox.ac.uk/fsl/fslwiki/

randomize) to perform permutation tests and generate a null distribution of cluster masses for multiple-comparisons correction (using FSL's "threshold-free cluster enhancement," *p* < 0.05 two tailed).

## Results

### Participants learned the animal likelihoods in the tours task
We evaluated participants' final learning of the likelihood of each animal in each sector using performance from the last set of tours on the last day. In those tours, the participants chose the more likely animal 73 ± 3% of the time. Note that even if participants had perfect knowledge of the animal likelihoods, we might not expect participants to choose the more likely animal 100% of the time, due to noise in the decision process or probability matching (the tendency to match choice probabilities to the probability of each option being correct; Vulkan, 2000; Erev and Barron, 2005).

In a postexperiment questionnaire, we asked participants to estimate the animal likelihoods *P(animal | sector)* for every animal and every sector. These estimates were close to the true likelihoods, on average (Fig. 4A). The mean Kullback–Leibler divergence of the estimated likelihoods from the real likelihoods was 0.13 ± 0.015. As discussed below, we used these participant-estimated likelihoods in our neural analyses, in lieu of the correct likelihoods.

### Performance on photographs task suggested maintenance of posterior distributions over sectors
During the fMRI scan sessions, participants correctly chose the more (or less) probable sector 67 ± 1% of the time, which is significantly above chance ($t_{(31)} = 13.1$, $p < $ 1e-12). Moreover, logistic regression on participants' responses showed that, the larger the difference in posterior probability between the correct and incorrect options, the more likely participants were to choose the correct answer (Fig. 4B). Again, as in the tours task, we expected stochasticity in participants' behavior due to noise and probability matching.

Note that the two sector options in each question were chosen at random, and therefore required participants to discriminate between posterior probabilities for any possible pair of sectors. Interestingly, participants performed similarly well
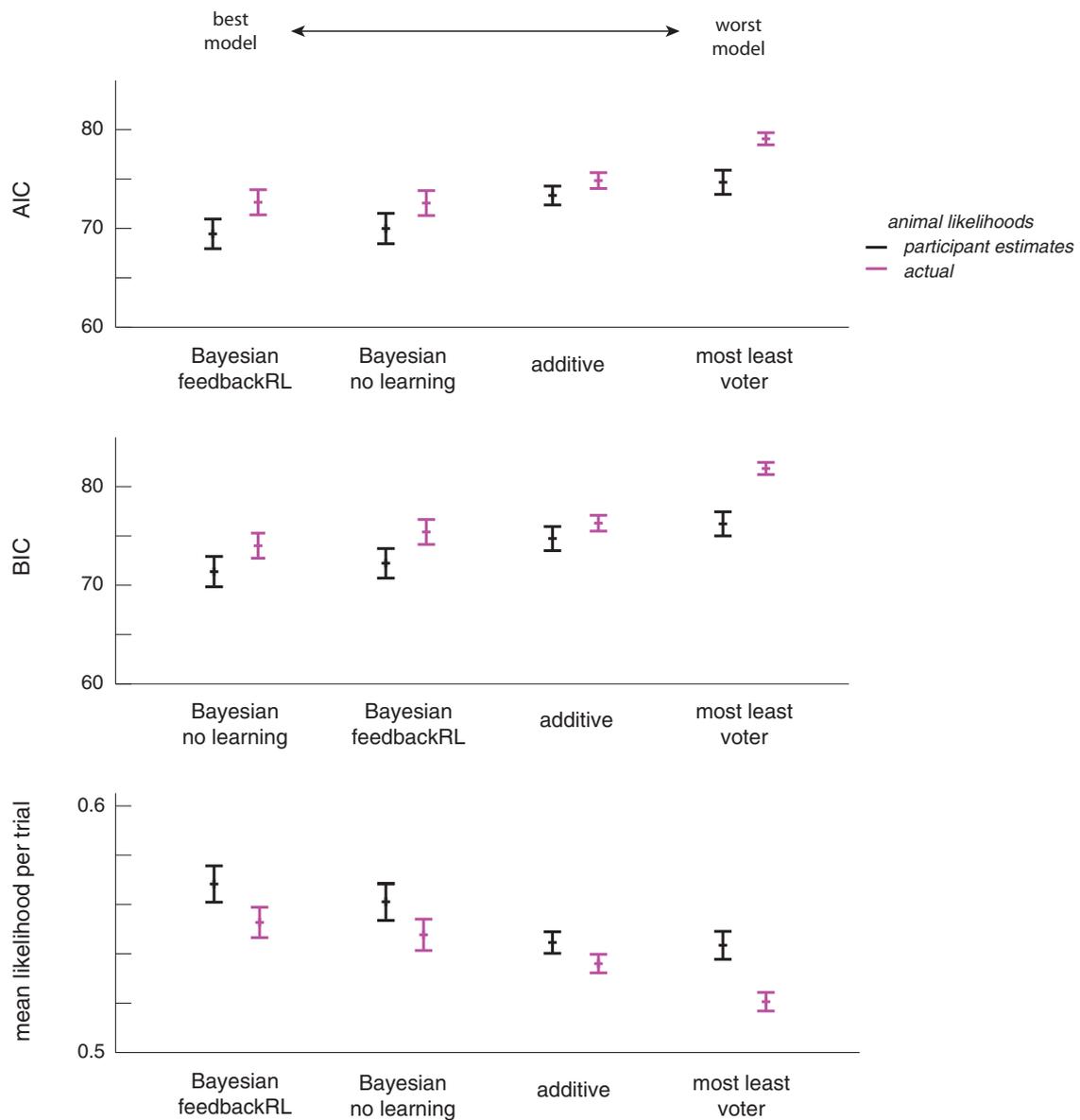
**Figure 5.** Behavioral model fitting. AIC, BIC, and (geometric) mean likelihood per trial (i.e., the exponentiated mean log likelihood per trial, without penalization for number of parameters) for the best-fitting behavioral model in each class (mean ± SEM across participants) suggest that the Bayesian models explained the behavioral data best. Note that better behavioral model fits are indicated by low AIC and BIC scores, but high mean likelihood. Results are shown for model fits using the participant estimates of the likelihoods or using the actual (true) likelihoods.

whether or not questions included the MAP (most probable) sector (accuracy 69 ± 2% for questions including the MAP; 66 ± 1% for questions not including the MAP; not significantly different). This result further indicates that participants were tracking the full posterior distribution, and not just the MAP sector.

**Trial-by-trial behavioral model fitting suggested that participants were approximately Bayesian**
The relative performance of the behavioral models is shown in Figure 5, and the mean parameter fits are shown in Table 2. For model comparison, we used the best-performing version from each class of behavioral models (Table 2).

The two Bayesian behavioral models (with and without feedbackRL) performed best, explaining the data about equally well. Overall, the model with feedbackRL was the best behavioral model according to AIC, but the Bayesian model without learn-

ing was the best behavioral model according to BIC, which penalizes more strongly for extra parameters.

The additive behavioral model performed worse than the Bayesian behavioral models, indicating that participants were accumulating evidence multiplicatively, in accordance with the optimal strategy (Eq. 1). None of the heuristic inference models that we tested (the "most–least voter" class of models) could successfully outperform the Bayesian behavioral models. Nor did we identify any significant effect of recency or primacy (any small improvements in the model likelihoods were not justified by the increased number of parameters). We therefore concluded that participants were Bayesian or near-Bayesian in their inference.

As shown in Figure 5, using the participants' subjective estimates of the animal likelihoods (from the postexperiment questionnaire) provided a better fit for all behavioral models, compared with using the real animal likelihoods. This may be surprising for the feedbackRL behavioral model, given that the

**Table 2. Free parameters and parameter fits, for the best-fitting behavioral model for each class[a]**

| Behavioral model | Free parameters | Mean $\pm$ SE | Range |
|---|---|---|---|
| Bayesian_nolearning | $\beta$: softmax inverse temperature | 4.04 $\pm$ 2.22 | [0, $\infty$] |
| Additive | $\beta$: softmax inverse temperature | 7.04 $\pm$ 3.46 | [0, $\infty$] |
| Mostleast_voter (voting for or against the sectors in which an animal was the most or least common) | $v\_pos$: size of positive vote; $v\_neg$: size of negative vote | 1.69 $\pm$ 3.39; 0.754 $\pm$ 1.39 | [0, $\infty$]; [0, $\infty$] |
| Bayesian_feedbackRL (learning from all feedback, no scaling of learning rates, and $\alpha_{pos} = \alpha_{neg}$) | $\alpha$: learning rate; $\beta$: softmax inverse temperature | 0.0515 $\pm$ 0.161; 4.82 $\pm$ 2.52 | [0, 1]; [0, $\infty$] |

[a]The best-fitting behavioral models for all classes did not model recency or primacy biases, and used each participant's subjective estimates of the animal likelihoods rather than the actual likelihoods. For model classes that had additional variants, the best-fitting settings are described in parentheses.

participant estimates were elicited at the end of the experiment, but were used in the model to initialize estimates of the likelihoods. However, the low learning rates (Table 2, average fit learning rates; also, 19% of participants had fitted learning rates of 0) suggest that changes in the likelihoods throughout the experiment were small relative to the differences between the real and estimated likelihoods. The low learning rates also explain why the feedbackRL behavioral model fit the data similarly well to a Bayesian behavioral model that did not allow for changes of the likelihood during the task: the models are nested (identical for learning rates of zero) and similar for low learning rates.

**Representational similarity analysis suggests that the suborbital sulcus contains a representation of the (log) posterior distribution over latent causes**

Figure 6B shows the representational similarity match of the suborbital sulcus with each of the neural models, relative to the representational similarity match with the best neural model—the logposterior (absolute values of the representational similarity match are shown in Table 4). For all of the alternative neural models tested, 95% or more of our bootstrap samples showed better representational similarity match for the logposterior than for the alternative neural model.

Because the posterior distribution tends to be more similar for neighboring time points compared with more distant time points, and that might also be the case for neural patterns, we took special care to verify that the logposterior neural model was superior to the alternative (control) time neural model. This was indeed the case. Moreover, we found that the temporal model displayed a negative representational similarity match with the neural patterns, because BOLD patterns for neighboring time points tended to be anticorrelated. This result was not dependent on our linear model for temporal distances; because we used Spearman's rank correlation to compute representational similarity match, the negative similarity match result would be observed for any other model of temporal distance that falls off monotonically (e.g., an exponential model). Therefore, since the posterior distribution showed a positive similarity match while the temporal neural model showed a negative similarity match, we can conclude that any positive correlations between the similarity matrices for the posterior distribution and time cannot be responsible for the representational similarity result for the posterior distribution.

Searchlight results for the representational similarity analysis are shown in Figure 7. The OFC and ventromedial prefrontal cortex showed a significantly greater representational similarity match for the logposterior model compared with every other neural model ($p < 0.05$ corrected, for every comparison), except for the entropy, posterior ranking, and posterior models. It also showed a greater representational similarity match for the log-

posterior than entropy using a more liberal threshold of $p < 0.05$ uncorrected.

## Discussion

Because the underlying structure of the world is often not directly observable, we must make inferences about the underlying situations or "latent causes" that generate our observations. The statistically optimal way to do this is to use Bayes' rule to infer the posterior distribution over latent causes. Based on previous studies implicating the OFC in the representation of the current context or situation (related to the idea of state in studies of reinforcement learning and decision making, and to the idea of schemas in studies of episodic memory), we hypothesized that the OFC might represent a posterior probability distribution over latent causes, computed using approximately Bayesian inference. To test this, we asked participants to make inferences about the probability of possible situations in an environment where the situation probabilistically generated their observations.

Using representational similarity analysis of fMRI activity during the inference task, we found that patterns of activity in the suborbital sulcus within the OFC were indeed best explained as representing a posterior distribution over latent causes. Searchlight analyses implicated the OFC more generally in this representation. Furthermore, participants' behavioral performance showed that they had access to a full posterior distribution over the latent causes for their choices; using trial-by-trial model fitting, we showed that participants' behavior was best explained as using Bayesian inference.

Our study provides evidence that the OFC represents a full posterior distribution over situations, as opposed to the best guess of the situation (the MAP) or other summary measures of the distribution, such as the overall uncertainty. We operationalized uncertainty as the entropy of the distribution; the highest entropy occurs when the distribution is completely flat (i.e., the participant is maximally uncertain about which latent cause generated the observations), and the lowest entropy occurs when the distribution is fully loaded on one latent cause (i.e., the participant is absolutely certain about which latent cause generated the observations). Our similarity analyses showed the entropy to have a widespread positive similarity match in many areas of the cortex, which we might expect because entropy should be correlated with the difficulty of the task, and so entropy might therefore be correlated with greater overall activity in many regions of the brain. Nonetheless, in >95% of our bootstrap samples, activity in the OFC was better explained by the posterior distribution than by the entropy. Furthermore, searchlight analyses showed the specificity of this result.

Our results, using multivariate analysis, build on previous fMRI studies that have used univariate analyses in the OFC to
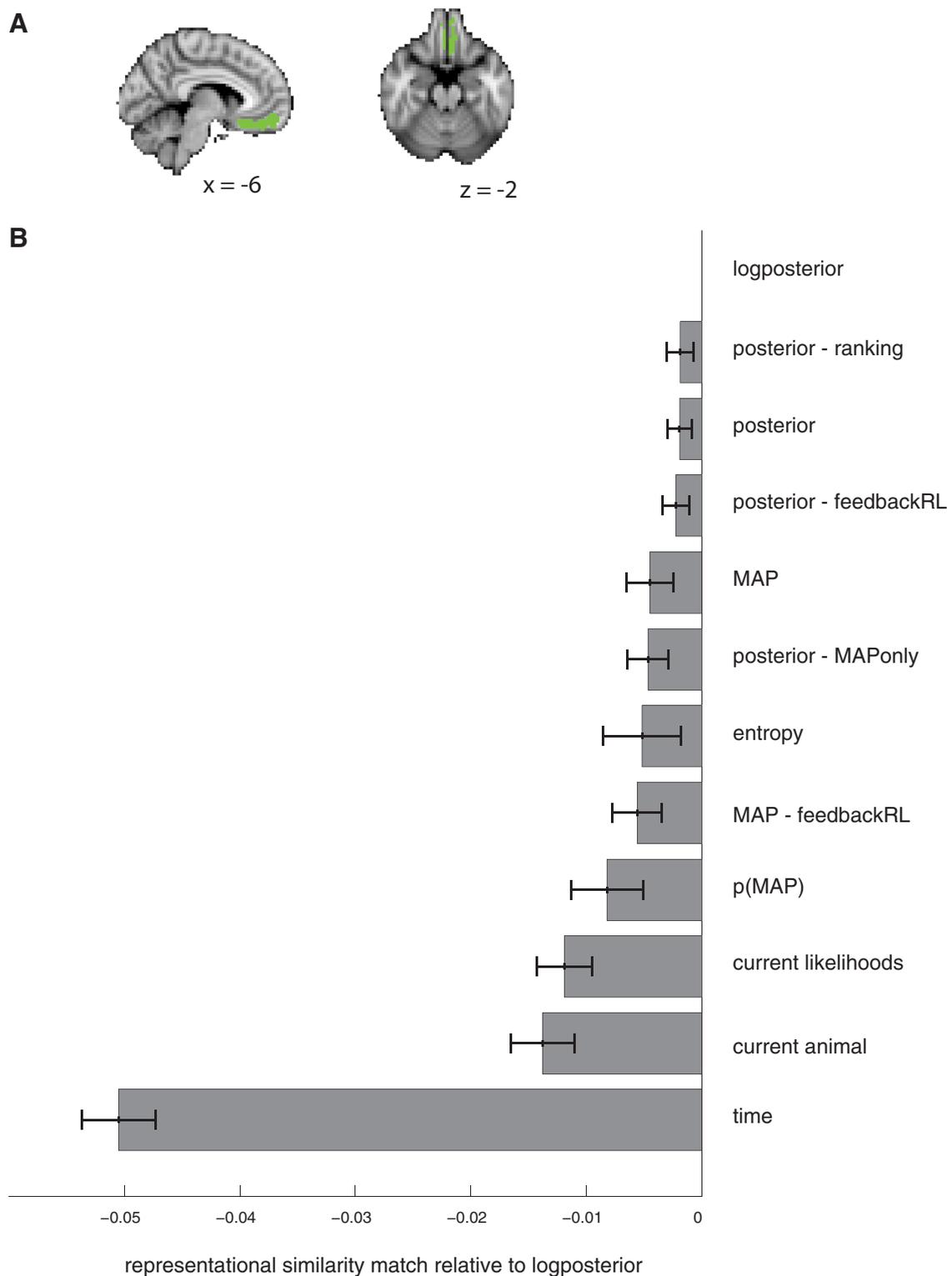
**A**



x = -6          z = -2

**B**



representational similarity match relative to logposterior

**Figure 6.** Representational similarity match for each neural model in the ROI. *A*, ROI: the suborbital cortex, a medial subregion of the OFC. See Materials and Methods for a description of how the region was defined. *B*, Representational similarity match for each of the neural models tested (Table 3), relative to the best neural model for the data (the logposterior), ordered by mean representational similarity match. The logposterior model showed the highest mean representational similarity match. The error bars indicate SEs of within-participant differences for each of the neural models compared with the logposterior. For all of the alternative neural models tested, ≥95% of our bootstrap samples showed a better match for the logposterior than for the alternative neural model.

investigate a range of summary statistical quantities related to the posterior distribution, but which do not capture the full distribution. These studies have shown that univariate activation of the ventromedial prefrontal cortex (vmPFC; which includes or is similar to our ROI) is correlated with a variety of summary sta-

tistics, e.g., expected reward (Ting et al., 2015), reward uncertainty (Critchley et al., 2001; Tobler et al., 2007), variance of the prior distribution in a sensory task (Vilares et al., 2012), and marginal likelihood of the current stimulus (d'Acremont et al., 2013). Our experiment benefited from several key features:

**Table 3. Neural models used in the representational similarity analysis, and the similarity measure used to derive the similarity matrix**

| Neural model | Description | Similarity measure for two timepoints |
|---|---|---|
| Posterior | Vector [4 × 1] containing the posterior probability of each sector, $P(sector \mid animals\ seen\ so\ far)$ | Normalized correlation* |
| Posterior–ranking** | Vector [4 × 1] containing the posterior probability of each sector, $P(sector \mid animals\ seen\ so\ far)$ | Spearman correlation |
| Log posterior | Vector [4 × 1] containing the natural logarithm of the posterior probability for each sector, $log[P(sector \mid animals\ seen\ so\ far)]$ | Normalized correlation* |
| Current animal | An integer $\in \{1, 2, 3, 4, 5\}$ indicating which animal is currently on screen | 1 if the same animal; 0 otherwise |
| Entropy | A scalar indicating the entropy of the posterior distribution over sectors | $-abs[entropy(t_1) - entropy(t_2)]$ |
| MAP | An integer $\in \{1, 2, 3, 4\}$ indicating which sector has the highest posterior probability | 1 if the same sector; 0 otherwise |
| $p$(MAP) | A scalar indicating the probability of the MAP sector | $-abs[p(MAP(t_1)) - p(MAP(t_2))]$ |
| Posterior–MAPonly | The posterior [4 × 1], zeroed for all sectors except the MAP sector (i.e. a signal that contains both MAP and $p$(MAP) information) | Normalized correlation* |
| Time | A scalar indicating the seconds passed since the start of the session | $-abs[time_1 - time_2]$ |
| Posterior–feedbackRL | Vector [4 × 1] indicating the posterior distribution over sectors, computed using the likelihoods updated on each trial using the best-fitting feedbackRL behavioral model (free parameters fitted for each participant) | Normalized correlation* |
| MAP–feedbackRL | An integer $\in \{1, 2, 3, 4\}$ indicating the most probable sector according to the best-fitting feedbackRL behavioral model (free parameters fitted to each participant) | 1 if the same sector; 0 otherwise |

*The normalized correlation of vectors x and y is $x \cdot y/(\|x\| * \|y\|)$, and is equivalent to the cosine of the angle between the two vectors. It behaves differently than the more commonly used Pearson correlation; for example, the posterior distributions [0.24 0.25 0.25 0.26] and [0.26 0.25 0.25 0.24] have Pearson correlation of −1 but normalized correlation of 0.9994. We used normalized correlation because this measure accords better with intuition regarding the similarity of posterior distributions and quantities derived from posterior distributions; however, similar results were observed when using Pearson correlations instead.

**This neural model is the same as the "posterior" model, except that, by using the Spearman correlation as its similarity metric, it only retains information about the rank ordering of the sectors in the posterior distribution. It is identical to a rank ordering of the sectors in the logposterior model, since the logarithm is a monotonic operator.

**Table 4. Representational similarity match for each neural model in the ROI, ordered by mean representational similarity match[a]**

| Neural model | Representational similarity match |
|---|---|
| Logposterior | $0.01410 \pm 0.00291$ |
| Posterior–ranking | $0.01223 \pm 0.00283$ |
| Posterior | $0.01219 \pm 0.00284$ |
| Posterior–feedbackRL | $0.01185 \pm 0.00288$ |
| MAP | $0.00961 \pm 0.00244$ |
| Posterior–MAPonly | $0.00944 \pm 0.00246$ |
| Entropy | $0.00892 \pm 0.00243$ |
| MAP–feedbackRL | $0.00851 \pm 0.00217$ |
| $p$(MAP) | $0.00589 \pm 0.00187$ |
| Current likelihoods | $0.00220 \pm 0.00157$ |
| Current animal | $0.00033 \pm 0.00154$ |
| Time | $-0.03639 \pm 0.00404$ |

[a]This table shows mean and SE of the absolute values of the representational similarity match, rather than the mean and SE of the within-participant differences relative to the best neural model (shown in Fig. 6).

(1) multivariate neural analysis, (2) four different latent causes, and (3) dissociation of latent cause from both reward and motor plan. These features enabled us to identify orbitofrontal representation of a full posterior distribution over latent causes that was separate from value, and which also explained neural activity in the area better than any single summary statistic we tried. Our result may therefore explain why evidence for different summary statistics was found in different studies: these are all components of the full posterior distribution or correlates of it.

Our study also builds on previous work in the fields of decision making and episodic memory that has implicated the OFC in representations of the current situation or context. In decision making, a belief distribution over states enables one to optimally learn or compute a behavioral policy even when the state of the world is not directly observable (as in partially observable Markov decision processes; Kaelbling et al., 1998). The OFC has long been implicated in reinforcement learning and decision making in a wide range of settings; a recent review provides a



**Figure 7.** Whole-brain searchlight result. Brain areas that passed both of the following criteria: (1) significantly higher representational similarity match with the logposterior neural model compared with every other neural model from Table 3 except the posterior, the posterior computed using the feedbackRL behavioral model, the posterior ranking, and the entropy, at $p < 0.05$ with whole-brain correction for every comparison; (2) higher representational similarity match with the logposterior compared with the entropy, at $p < 0.05$ uncorrected. The map is displayed on the orbital/ventral surface of an inflated brain.

unifying explanation for these results by postulating that the OFC represents inferred states in partially observable situations (Wilson et al., 2014). In theories of episodic memory, it is believed that we organize our memories according to an inferred schema that specifies the situation and stores previously learned relationships that a new memory can be incorporated into (Tse et al., 2007; Hupbach et al., 2008). These schemas seem to be represented or processed in the vmPFC (Ranganath and Ritchey, 2012; van Kesteren et al., 2012; for review, see Schlichting and Preston, 2015), an area of the brain that is similar to our ROI. For example, Tse et al. (2011) showed evidence that activation of rat mPFC is highest immediately after memory encoding that should involve incorporating new information into existing schemas. Ezzyat and Davachi (2011) showed that greater activation of vmPFC in humans during memory encoding is correlated with how strongly those memories are associated with other memories in the same "event," which is consistent with the idea that the vmPFC is involved in schemas that are bound to memories. Our results confirm the involvement of the OFC in representations of the current situation, and additionally show that this representation in the OFC takes the form of a distribution over possible situations.

Finally, our work also builds on previous studies investigating neural circuits involved in the "weather prediction" task, very similar to ours, in which one of two "weather" outcomes is probabilistically predicted by sequences of cards. Knowlton et al. (1996) implicated the striatum in the learning of these probabilistic associations. In our task, participants learned the animal likelihoods outside the MR scanner, and thus we could not assess the brain areas involved in the learning phase. However, our results are compatible with those of Knowlton et al. insofar as the OFC may use associations learned by the striatum (in our experiment, the animal likelihoods) to make inferences when presented with new observations (in our experiment, the photographs task). More recently, Yang and Shadlen (2007) used the weather-prediction task to show representation of a decision variable in the parietal cortex that took the form of the log likelihood ratio between two options. In our experiment, we decorrelated the posterior probability from both decision variables and stimulus–reward associations, and we also investigated representations of the posterior probability over latent causes before the decision period. We conjecture that the OFC contains representations of the current state or situation in terms of a posterior distribution over the possible states, a representation that is likely used by downstream areas, e.g., parietal cortex, for decision making.

Previous work on the weather-prediction task also showed that most individuals used heuristic strategies in inferring the weather (Gluck et al., 2002). In our experiment, we explored several heuristic behavioral models of participants' inference, but were not able to find any that predicted participants' behavior better than the optimal Bayesian models. There are several reasons why our task may have discouraged the use of heuristics. First, the animal likelihoods in our experiment were designed to avoid one-to-one mappings between observations and latent causes. Second, the task environment had four possible latent causes (instead of two), and the task itself required rank ordering all four latent causes rather than just estimating the MAP, thus increasing complexity and leading to the inadequacy of simple heuristics. Finally, we provided participants with extensive training on the probabilistic relationships in the safari, so that heuristics may have been less necessary.

Neurally, the posterior distribution we found in the OFC was best modeled as being represented in log space. Representation in log space may be advantageous because addition can then replace the multiplicative operation required to accumulate evidence in nonlog space (e.g., across animal presentations, in our experiment); the ability of neurons to add is well characterized, while it is less clear to what extent neurons can multiply (Yuste and Tank, 1996; Peña and Konishi, 2001; Gabbiani et al., 2002). Indeed, neural representation in log space is common in many domains, e.g., decision variables (Yang and Shadlen, 2007), time (Gibbon, 1977), and numbers (Longo and Lourenco, 2007).

To summarize, we designed a task in which participants' observations were probabilistically generated by unobserved "situations" or "latent causes," and found evidence that the OFC represents a probability distribution over possible latent causes. A representation of the log posterior distribution explained OFC activity better than alternatives such as the best guess of the current situation or overall uncertainty in the current situation. This finding was further supported by behavioral evidence that participants had access to the full probability distribution for decision making, and used Bayesian inference to compute the probability distribution. Our results may explain why previous studies of the OFC have found evidence for representation of various summary statistical quantities in the OFC (these are in fact components of the full posterior probability distribution). Our results may also unify findings from disparate literature on reinforcement learning and episodic memory, which separately implicate the OFC in representations of the current situation.

## References

Courville AC, Gordon GJ, Touretzky DS, Daw ND (2004) Model uncertainty in classical conditioning. In: Advances in neural information processing systems 16. pp 977–984. Cambridge, MA: MIT.

Critchley HD, Mathias CJ, Dolan RJ (2001) Neural activity in the human brain relating to uncertainty and arousal during anticipation. Neuron 29:537–545. CrossRef Medline

d'Acremont M, Fornari E, Bossaerts P (2013) Activity in inferior parietal and medial prefrontal cortex signals the accumulation of evidence in a probability learning task. PLoS Comput Biol 9:e1002895. CrossRef Medline

Deichmann R, Gottfried JA, Hutton C, Turner R (2003) Optimized EPI for fMRI studies of the orbitofrontal cortex. Neuroimage 19:430–441. CrossRef Medline

Destrieux C, Fischl B, Dale A, Halgren E (2010) Automatic parcellation of human cortical gyri and sulci using standard anatomical nomenclature. Neuroimage 53:1–15. CrossRef Medline

Efron B, Tibshirani R (1986) Bootstrap materials and methods for standard errors, confidence intervals, and other measures of statistical accuracy. Stat Sci 1:54–75.

Erev I, Barron G (2005) On Adaptation, Maximization, and Reinforcement Learning Among Cognitive Strategies. Psychological Rev 112:912–931. CrossRef

Ezzyat Y, Davachi L (2011) What constitutes an episode in episodic memory? Psychol Sci 22:243–252. CrossRef Medline

Gabbiani F, Krapp HG, Koch C, Laurent G (2002) Multiplicative computation in a visual neuron sensitive to looming. Nature 420:320–324. CrossRef Medline

Gershman SJ, Niv Y (2010) Learning latent structure: carving nature at its joints. Curr Opin Neurobiol 20:251–256. CrossRef Medline

Gershman SJ, Blei DM, Niv Y (2010) Context, learning, and extinction. Psychol Rev 117:197–209. CrossRef Medline

Ghosh VE, Gilboa A (2014) What is a memory schema? A historical perspective on current neuroscience literature. Neuropsychologia 53:104–114. CrossRef Medline

Gibbon J (1977) Scalar expectancy theory and Weber's law in animal timing. Psychol Rev 84:279–325. CrossRef

Gluck MA, Shohamy D, Myers C (2002) How do people solve the "weather prediction" task?: individual variability in strategies for probabilistic category learning. Learn Mem 9:408–418. CrossRef Medline

Hupbach A, Hardt O, Gomez R, Nadel L (2008) The dynamics of memory: context-dependent updating. Learn Mem 15:574–579. CrossRef Medline

Kaelbling LP, Littman ML, Cassandra AR (1998) Planning and acting in partially observable stochastic domains. Artif Intell 101:99–134. CrossRef

Knowlton BJ, Mangels JA, Squire LR (1996) A neostriatal habit learning system in humans. Science 273:1399–1402. CrossRef Medline

Kriegeskorte N, Mur M, Bandettini P (2008) Representational similarity analysis-connecting the branches of systems neuroscience. Front Sys Neurosci 2:4. CrossRef Medline

Libby LA, Ekstrom AD, Ragland JD, Ranganath C (2012) Differential connectivity of perirhinal and parahippocampal cortices within human hippocampal subregions revealed by high-resolution functional imaging. J Neurosci 32:6550–6560. CrossRef Medline

Longo MR, Lourenco SF (2007) Spatial attention and the mental number line: Evidence for characteristic biases and compression. Neuropsychologia 45:1400-1407.

Peña JL, Konishi M (2001) Auditory spatial receptive fields created by multiplication. Science 292:249–252.

Ranganath C, Ritchey M (2012) Two cortical systems for memory-guided behaviour. Nat Rev Neurosci 13:713–726. CrossRef Medline

Richards BA, Xia F, Santoro A, Husse J, Woodin MA, Josselyn SA, Frankland PW (2014) Patterns across multiple memories are identified over time. Nat Neurosci 17:981–986. CrossRef Medline

Schlichting ML, Preston AR (2015) Memory integration: neural mechanisms and implications for behavior. Curr Opin Behav Sci 1:1–8. CrossRef Medline

Ting CC, Yu CC, Maloney LT, Wu SW (2015) Neural mechanisms for integrating prior knowledge and likelihood in value-based probabilistic inference. J Neurosci 35:1792–1805. CrossRef Medline

Tobler PN, O'Doherty JP, Dolan RJ, Schultz W (2007) Reward value coding distinct from risk attitude-related uncertainty coding in human reward systems. J Neurophysiol 97:1621–1632. Medline

Tse D, Langston RF, Kakeyama M, Bethus I, Spooner PA, Wood ER, Witter MP, Morris RG (2007) Schemas and memory consolidation. Science 316:76–82. CrossRef Medline

Tse D, Takeuchi T, Kakeyama M, Kajii Y, Okuno H, Tohyama C, Bito H, Morris RG (2011) Schema-dependent gene activation and memory encoding in neocortex. Science 333:891–895. CrossRef Medline

van Kesteren MT, Ruiter DJ, Fernández G, Henson RN (2012) How schema and novelty augment memory formation. Trends Neurosci 35:211–219. CrossRef Medline

Vilares I, Howard JD, Fernandes HL, Gottfried JA, Kording KP (2012) Differential representations of prior and likelihood uncertainty in the human brain. Curr Biol 22:1641–1648. CrossRef Medline

Vulkan N (2000) An economist's perspective on probability matching. J Econ Surv 14:101–118. CrossRef

Wilson RC, Takahashi YK, Schoenbaum G, Niv Y (2014) Orbitofrontal cortex as a cognitive map of task space. Neuron 81:267–279. CrossRef Medline

Yang T, Shadlen MN (2007) Probabilistic reasoning by neurons. Nature 447:1075–1080. CrossRef Medline

Yuste R, Tank DW (1996) Dendritic integration in mammalian neurons a century after Cajal. Neuron 16:701–716. CrossRef Medline