

## Brain-wide representations of prior information in mouse decision-making

Charles Findling<sup>1,19</sup>, Felix Hubert<sup>1</sup>, International Brain Laboratory\*, Luigi Acerbi<sup>2</sup>, Brandon Benson<sup>3</sup>, Julius Benson<sup>4</sup>, Daniel Birman<sup>5</sup>, Niccolò Bonacchi<sup>6</sup>, Matteo Carandini<sup>7</sup>, Joana A Catarino<sup>6</sup>, Gaëlle A Chapuis<sup>1</sup>, Anne K Churchland<sup>8</sup>, Yang Dan<sup>9</sup>, Eric EJ DeWitt<sup>6</sup>, Tatiana A Engel<sup>10</sup>, Michele Fabbri<sup>6</sup>, Mayo Faulkner<sup>7</sup>, Ila Rani Fiete<sup>11</sup>, Laura Freitas-Silva<sup>6</sup>, Berk Gerçek<sup>1</sup>, Kenneth D Harris<sup>7</sup>, Michael Häusser<sup>7</sup>, Sonja B Hofer<sup>12</sup>, Fei Hu<sup>9</sup>, Julia M Huntenburg<sup>13</sup>, Anup Khanal<sup>8</sup>, Chris Krasniak<sup>14</sup>, Christopher Langdon<sup>10</sup>, Peter E Latham<sup>15</sup>, Petrina Y P Lau<sup>7</sup>, Zach Mainen<sup>6</sup>, Guido T Meijer<sup>6</sup>; Nathaniel J Miska<sup>12</sup>, Thomas D Mrsic-Flogel<sup>12</sup>, Jean-Paul Noel<sup>4</sup>, Kai Nylund<sup>5</sup>, Alejandro Pan-Vazquez<sup>10</sup>, Liam Paninski<sup>16</sup>, Jonathan Pillow<sup>10</sup>, Cyrille Rossant<sup>7</sup>, Noam Roth<sup>5</sup>, Rylan Schaeffer<sup>11</sup>, Michael Schartner<sup>6</sup>, Yanliang Shi<sup>10</sup>, Karolina Z Socha<sup>7</sup>, Nicholas A Steinmetz<sup>5</sup>, Karel Svoboda<sup>17</sup>, Charline Tessereau<sup>13</sup>, Anne E Urai<sup>18</sup>, Miles J Wells<sup>7</sup>, Steven Jon West<sup>12</sup>, Matthew R Whiteway<sup>16</sup>, Olivier Winter<sup>6</sup>, and Ilana B Witten<sup>10</sup>, Anthony Zador<sup>14</sup>, Peter Dayan<sup>13</sup>, Alexandre Pouget<sup>1</sup>.

<sup>1</sup>University of Geneva, Switzerland; \*Virtual Entity; <sup>2</sup>University of Helsinki, <sup>3</sup>Stanford University, USA; <sup>4</sup>New York University, USA; <sup>5</sup>University of Washington, USA; <sup>6</sup>Champalimaud Foundation, Portugal; <sup>7</sup>University College London, UK; <sup>8</sup>University of California Los Angeles, USA; <sup>9</sup>University of California Berkeley, USA; <sup>10</sup>Princeton University, USA; <sup>11</sup>Massachusetts Institute of Technology, USA; <sup>12</sup>Sainsbury Wellcome Center, University College London, UK; <sup>13</sup>Max Planck Institute, University of Tübingen, Germany; <sup>14</sup>Cold Spring Harbor Laboratory; <sup>15</sup>Gatsby Computational Neuroscience Unit, UK; <sup>16</sup>Columbia University, USA; <sup>17</sup>Allen Institute for Neural Dynamics, USA; <sup>18</sup>Leiden University, The Netherlands

<sup>19</sup> Correspondence: [charles.findling@internationalbrainlab.org](mailto:charles.findling@internationalbrainlab.org)

**The neural representations of prior information about the state of the world are poorly understood. To investigate this issue, we examined brain-wide Neuropixels recordings and widefield calcium imaging collected by the International Brain Laboratory. Mice were trained to indicate the location of a visual grating stimulus, which appeared on the left or right with prior probability alternating between 0.2 and 0.8 in blocks of variable length. We found that mice estimate this prior probability and thereby improve their decision accuracy. Furthermore, we report that this subjective prior is encoded in at least 20% to 30% of brain regions which, remarkably, span all levels of processing, from early sensory areas (LGd, VISp) to motor regions (MOs, MOp, GRN) and high level cortical regions (ACCd, ORBvl). This widespread representation of the prior is consistent with a neural model of Bayesian inference involving loops between areas, as opposed to a model in which the prior is incorporated only in decision making areas. This study offers the first brain-wide perspective on prior encoding at cellular resolution, underscoring the importance of using large scale recordings on a single standardized task.**

The ability to combine sensory information with prior knowledge through probabilistic inference is crucial for perception and cognition. In simple cases, inference is performed near-optimally by the brain, following key precepts of Bayesian decision theory (Ernst & Banks, 2002; Jacobs, 1999; Knill & Pouget, 2004; Mamassian et al., 1998; Weiss et al., 2002). For example, when interpreting a visual scene, we assume *a priori* that light comes from above – a sensible assumption which allows us to resolve otherwise ambiguous images (Mamassian et al., 1998).

While much theoretical work has been devoted to the neural representation of Bayesian inference (Echeveste et al., 2020; Ganguli & Simoncelli, 2014; Ma et al., 2006; Soltani & Wang, 2010), it remains unclear where and how prior knowledge is represented in the brain. At one extreme, the brain might combine prior information with sensory evidence in high level decision-making brain regions, right before decisions are turned into actions. This would predict that prior information is encoded only in late stages of processing, as has indeed been reported in parietal, orbitofrontal and prefrontal cortical areas (Forstmann, 2010; Hanks et al., 2011; Hansen et al., 2012; Mulder et al., 2012; Niv, 2019; Nogueira et al., 2017; Rao et al., 2012). At the other extreme, the brain might operate like a very large Bayesian network, in which probabilistic inference is the *modus operandi* in all brain regions and inference can be performed in all directions (Ackley et al., 1985; Berkes et al., 2011; Bondy et al., 2018; Jardri et al., 2017; Kok et al., 2012; Lange & Haefner, 2022). This would allow neural circuits to infer beliefs over variables from observations of arbitrary combinations of other variables. For instance, upon seeing an object, the brain might be able to infer its auditory and tactile properties; but could just as well perform the reverse inference, i.e., predicting its visual appearance upon hearing or touching it. Such a model would predict that prior information should be available throughout the brain, even in low level cortical sensory areas (Berkes et al., 2011; Bondy et al., 2018; Kok et al., 2012; Lange & Haefner, 2022). The current literature offers a contradictory, and thus inconclusive, perspective on whether the prior is indeed encoded in brain regions associated with early processing (Bell et al., 2016; Bondy et al., 2018; Haefner et al., 2016; Han & Helmchen, 2023; Hanks et al., 2011; Ishizu et al., 2023; Mayrhofer et al., 2019; Park et al., 2022; Platt & Glimcher, 1999; Rao et al., 2012; Son et al., 2023). This is because past studies have collectively recorded from only a limited set of areas; and, since they use different tasks, even these results cannot be fully integrated.

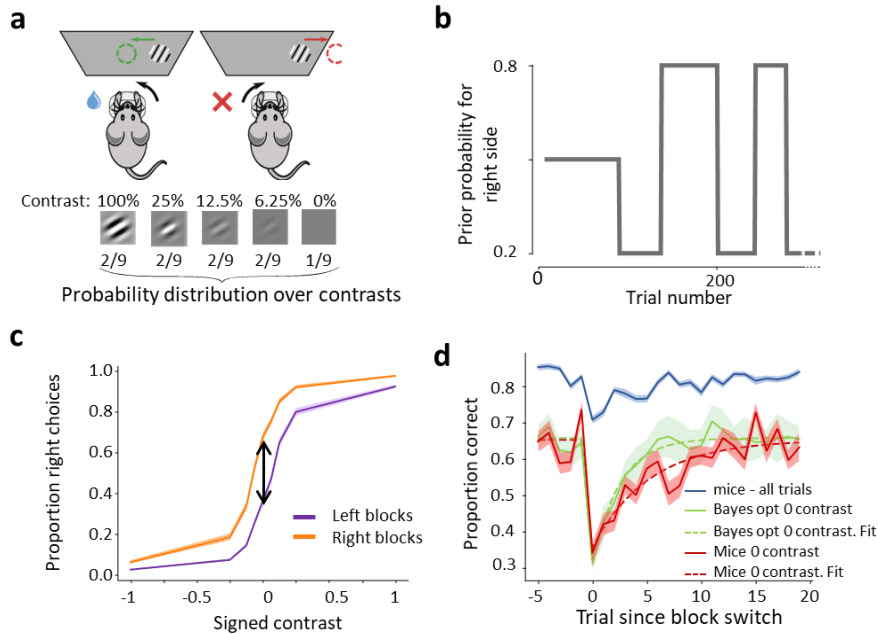
To address this problem, we analyzed brain-wide data from the International Brain Laboratory, which provides electrophysiological recordings from 159 brain regions, defined by the Allen Common Coordinate Framework (Wang et al., 2020), as well as from widefield imaging (WFI) data from layers 2/3 of cortex of mice performing the same decision-making task (International Brain Lab et al., 2023; The International Brain Laboratory et al., 2021). Our results suggest that the prior is encoded cortically and subcortically, across all levels of the brain, including early sensory regions.

### Mice use the prior to optimize their performance

Mice were trained to discriminate whether a visual stimulus, of varying contrast, appeared in the right or left visual field (Fig. 1a). Importantly, the prior probability that the stimulus appeared on the right side switched in a random and uncued manner between 0.2 and 0.8 in blocks of 20-100 trials (Fig. 1b). Knowledge of the current prior would help the mice perform well; in particular, on zero contrast trials, the prior is the only source of information, as the probability of reward on these trials is determined by the block probability. We refer to the experimentally-determined prior as the ‘true block prior’. Since the presence of the blocks was not explicitly cued, mice could only form a subjective estimate of the true block prior from the trial history. At best, they can compute the estimate of the true block prior given full knowledge of the task structure and the sequence of previous stimulus sides since the start of a session, which we refer thereafter as the Bayes-optimal prior (see Methods, Fig. 2a).

Analyzing choice behavior revealed that mice leverage the block structure to improve their performance. Psychometric curves conditioned on right and left blocks, averaged across all animals and all sessions, were displaced relative to each other, in a direction consistent with the true block prior (2-tailed signed-rank Wilcoxon paired test between proportion of right choices on zero contrast trials:  $t=3$ ,  $p=1.4E-20$ ,  $N=115$  mice; Fig. 1c). As a result, animals performed at  $59.0\% \pm 0.4\%$  (mean  $\pm$  sem) correct for zero contrast trials, statistically significantly better than chance (2-tailed signed-rank Wilcoxon  $t=12$ ,  $p=2.6E-20$ ,  $N=115$  mice) but significantly worse than an observer that generates actions by sampling from the Bayes-optimal prior, which performs at  $61.3\% \pm 1.8\%$  (mean  $\pm$  std; 2-tailed signed-rank Wilcoxon paired test  $t=1547$ ,  $p=6.0e-7$ ,  $N=115$  mice).

Tracking performance around block switches provided further evidence that the animals estimated and used the prior. Indeed, around block switches performance dropped, presumably because of the mismatch between the subjective and true block prior. Thus, performance on zero contrast trials recovered with a decay constant of 5.16 trials (jackknife median, see Methods). This is slower than the previously introduced observer that generates actions by sampling the Bayes-optimal prior (jackknife median: 2.46, 2-tailed paired t-test  $t_{114}=2.94$ ,  $p=0.004$ ,  $N=115$  jackknife replicates; Fig. S1).



**Figure 1:** Mice use the block prior to improve their performance. **a.** Mice had to move a visual grating, appearing  $35^\circ$  in the periphery (here shown on the right hand side), to the center of the screen by tuning a wheel with their front paws. The contrast of the visual stimulus varied from trial to trial. **b.** The prior probability that the stimulus appeared on the right side was maintained at either 0.2 or 0.8 over blocks of trials, after an initial block of 90 trials during which the prior was set to 0.5. The length of a block was drawn from a truncated exponential distribution between 20 and 100 trials, with the scale parameter of the exponential set to 60 trials. Following a wheel turn, the mouse was provided with positive feedback in the form of a water reward, or negative feedback in the form of a white noise pulse and timeout. The next trial began after a delay, followed by a quiescence period uniformly sampled between 400ms and 700ms during which the mice had to hold the wheel still. **c.** Psychometric curves averaged across all 115 animals and 354 sessions and conditioned on block identity. The proportion of right choices on zero contrast trials is significantly different across blocks (2-tailed Wilcoxon paired test:  $t=3$ ,  $p=1.4E-20$ ,  $N=115$ ) and displaced in the direction predicted by the true block prior (black double arrow). **d.** Reversal curves showing the percentage of correct responses following block switches. Blue: average performance across all animals and all contrasts. Dark red: same as pink but for zero contrast trials only. Green: performance of an observer generating choices stochastically according to the Bayes-optimal estimate of the prior on zero contrast trials. Shaded region around the mean shows the s.e.m across mice for the curves showing mouse behavior (blue and red curves) and the standard deviation for the Bayes-optimal model (green curve), as there is no inter-individual variability to account for.

### Decoding the prior during the inter-trial interval

In order to determine where the prior is encoded in the mouse's brain, we used linear regression to decode the Bayes-optimal prior from neural activity during the inter-trial interval (ITI) when wheel movements are minimized (from -600ms to -100ms before stimulus onset, see Methods, Fig. S2)(The International Brain Laboratory et al., 2021). Note that we do not decode the true block prior, since mice are not explicitly cued about block identity and, therefore cannot possibly know this quantity. We assess the quality of the decoding with an  $R^2$  measure. However, to assess the statistical significance of this value we cannot use standard linear regression methods, since these assume independence of trials,

while both neural activity and the prior exhibit temporal correlations. Instead, we use a “pseudosession” method (Harris, 2020): we first construct a null distribution by decoding the (counterfactual) Bayes-optimal priors computed from stimulus sequences generated by sampling from the same process as that used to generate the stimulus sequence that was actually shown to the mouse (see Methods). A session was deemed to encode the prior significantly if  $R^2$  computed for the actual stimuli is larger than the 95<sup>th</sup> percentile of the null distribution generated from pseudosessions; effect sizes are reported as a corrected  $R^2$ , the difference between the actual  $R^2$  and the median  $R^2$  of the null distribution. Using corrected  $R^2$  is important because any slow drift in the recordings, for instance from movements of the probes across trials, could yield a non-zero uncorrected  $R^2$ . All values of  $R^2$  reported in this paper are corrected  $R^2$  unless specified otherwise. Later on, we will assess significance after applying the Benjamini-Hochberg correction for multiple comparisons with a conservative false discovery rate of 1%.

For completeness, we decoded both the Bayes-optimal prior and its log odds ratio ( $\log(\hat{p}/(1 - \hat{p}))$  where  $\hat{p}$  is the Bayes-optimal prior for the right side). For the Bayes-optimal prior, the analysis of the electrophysiological data (Ephys) revealed that around 30% of brain areas, spanning forebrain, midbrain, and hindbrain, encoded the prior significantly (Fig. 2). For example, we could decode the Bayes-optimal prior from a population of 107 neurons in ventrolateral orbitofrontal cortex (ORBvl) with accuracy of  $R^2 = 0.16$  (Fig. 2a,  $p=0.005$ ). Overall we could decode the Bayes-optimal prior in 29.5% (47/159) of the brain regions ( $p<.05$ , pseudosession test; Fisher’s method to combine p-values from multiple recordings of one region, no multiple comparisons correction; Fig. 2b). These regions include associative cortical areas like the ORBvl and the dorsal anterior cingulate area (ACAd), as well as early sensory areas such as the primary visual cortex (VISp) or the lateral geniculate nucleus (LGd). The Bayes-optimal prior can also be decoded from cortical and subcortical motor areas, such as primary and secondary motor cortex, the intermediate layer of the superior colliculus (SCm), the gigantocellular reticular nucleus (GRN) and the pontine reticular nucleus (PRNr), even though we decoded activity during the inter-trial interval, when wheel movements are minimal (Fig. S2). The encoding of the Bayes-optimal prior is also visible in the PSTH of single neurons (Fig. S3). Similar results were obtained when decoding the log odds ratio of the Bayes-optimal prior, instead of the linear Bayes-optimal prior: 31.4% (50/159) of the decoded regions were found to encode the log odds ratio significantly, again at all levels of brain processing (Fig. S4).

The analysis of widefield calcium imaging data (WFI) suggests an even more widespread encoding of the prior in cortex. Indeed, the Bayes-optimal prior was found to be significantly reflected in all dorsal cortical regions (Fig. 2c). This result may reflect a better signal-to-noise ratio, but it might also be due to the calcium signal from axons arising outside these specific areas. However, we also found that the corrected region-specific  $R^2$  for WFI and Ephys modalities were significantly correlated (Spearman correlation  $R=0.57$ ,  $p=0.0007$ ,  $N=31$  regions - Fig. 2d) - even when correcting for region sizes (Fig. S5) - thus suggesting that the signals we decode are, at least partly, specific to the decoded regions.

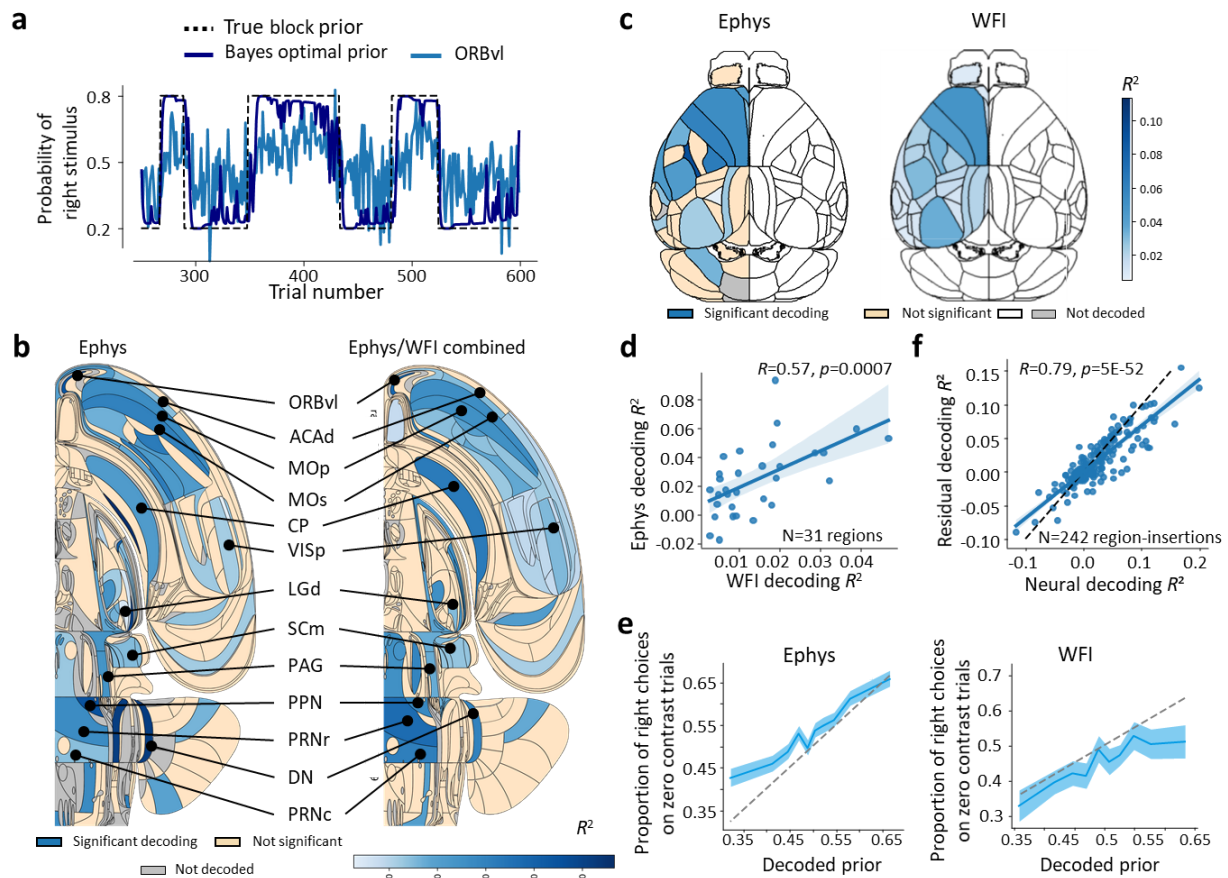
We also analyzed a larger data set containing all the units in the previous data set, plus units that had been identified through spike sorting but did not pass quality control criteria (QC=0, see Methods). This data set provides a near complete coverage of the mouse brain (267 regions instead of 159)(International Brain Lab et al., 2023). Although it contains less reliable spike trains, we once again observed that 30.7% of the regions (82/267) encode the Bayes-optimal prior significantly, distributed as before in sensory, motor and associative regions (Fig. S6a). The corrected  $R^2$  values were also highly correlated with the results obtained when using only the quality controlled units (Spearman correlation

$R=0.57$ ,  $p=4E-15$ ,  $N=159$  regions Fig. S6b). A smaller fraction of regions (21.3%, 57/267), but still at all levels of brain processing, were found to be significant when merging this larger Ephys data set and WFI data into a single map (using Fisher's method to combine p-values across Ephys and WFI, see Methods) and applying the Benjamini-Hochberg correction for multiple comparisons with a false discovery rate of 1%, Fig. 2b.

If the decoded prior is truly related to the subjective prior inferred and used by the animal, the amplitude of the decoded prior should be correlated with the animals' performance on zero contrast trials. Fig. 2e shows that this is indeed the case for both the Ephys and WFI data: on zero contrast trials, the probability that the mice chose the right side was proportional to the decoded Bayes-optimal prior (see Methods). Importantly, this relationship remained significant even after controlling for possible drift in the recordings (Fig. S7).

Our results indicate that the Bayes-optimal prior is encoded in multiple areas throughout the brain. However, it is conceivable that mice adjust their body posture or movement according to the subjective prior and that neural activity in some areas simply reflects these body adjustments. We call this an embodied prior. To test for this possibility, we analyzed video recordings, using Deep Lab Cut (DLC)(Laboratory, International Brain, 2022; Mathis et al., 2018) to estimate the position of multiple body parts, whisking motion energy, and licking during the inter-trial interval (ITI). We then trained a decoder of the Bayes-optimal prior from these features, and found significant decoding in 40.5% (43/106) of sessions. For these sessions, we found that the  $R^2$  for the prior decoded from video features was correlated with the  $R^2$  for the prior decoded from neural activity (at the brain region level), thus suggesting that the prior signal might be an embodied prior related to body posture (Spearman correlation  $R=0.13$ ,  $p=0.048$ ,  $N=242$  region insertions, Supp Fig. S8a). To test for this possibility further, we decoded the prior residual, defined as the Bayes-optimal prior minus the Bayes-optimal prior estimated from video features, from neural activity (again, at the brain region level). If the neural prior simply reflects the embodiment of features extracted by DLC, we should not be able to decode the prior residual from the neural activity and the  $R^2$  of the prior residual should not be correlated with the  $R^2$  of the full prior decoded from neural activity. Crucially, this is not what we observed. Instead, these two values of  $R^2$  are strongly correlated (Fig. 2f, Spearman correlation  $R=0.79$ ,  $p=5E-52$ ,  $N=242$  region insertions) thus suggesting that the neural prior is not an embodied prior, or at least that it cannot be fully explained by the motor features extracted from the video.

We also checked whether changes in eye position across blocks could account for the significant results in early visual areas such as VISp or LGd. It is indeed conceivable that mice look in the direction of the expected stimulus prior to a trial. If so, what we interpret as a prior signal might simply be due to a signal related to eye position. Consistent with this possibility, we found a significant correlation (Spearman correlation  $R=0.48$ ,  $p=0.024$ ,  $N=22$  region insertions) between the neural decoding  $R^2$  and the eye position decoding  $R^2$ , *i.e.*, the  $R^2$  for decoding the Bayes-optimal prior from eye position (using sessions in which video was available and recordings were performed in VISp and LGd,  $N = 22$ , Fig. S8b). Following the same approach as for the body posture and motion features, we then decoded the prior residual (Bayes-optimal prior minus Bayes-optimal prior estimated from eye position) from neural activity and found that the residual decoding  $R^2$  was correlated with the neural decoding  $R^2$  (Pearson correlation  $R=0.85$ ,  $p=7E-7$ ,  $N=22$  region insertions Fig. S8c). Therefore, the prior signals found in VISp and LGd are not simply reflecting subtle changes in eye position across blocks.



**Figure 2:** Encoding of the prior across the brain during the inter-trial interval. **a.** Bayes-optimal prior (dark blue) versus Bayes-optimal prior decoded from the ventrolateral orbitofrontal (ORBvl, light blue) on one specific session (corrected  $R^2=0.16$ , uncorrected  $R^2=0.25$ ). Dashed black line: true block prior. **b.** Left: Swanson map of cross-validated  $R^2$  for Ephys data. Right:  $R^2$  averaged across Ephys and WFI for areas that have been deemed significant (using Fisher’s method for combining  $p$ -values, see Methods). The map on the right includes a larger Ephys data set containing all units regardless of whether they pass quality control criteria (see Methods). Significance is assessed with the Benjamini-Hochberg procedure, correcting for multiple comparisons, with a conservative false discovery rate of 1%. 29.5% (left) and 21.3% (right) of the areas encode the prior significantly, at all levels of brain processing in both cases. For the full names of brain regions see [online table](#). **c.** Comparison of Ephys and WFI results for dorsal cortex. All areas significantly encode the Bayes-optimal prior in the WFI data. **d.** The corrected  $R^2$  for Ephys and WFI are significantly correlated (Spearman correlation  $R=0.57$ ,  $p<0.001$ ,  $N=31$  regions). Each dot corresponds to one region of the dorsal cortex. Significant and non-significant Ephys  $R^2$  were included in this analysis. **e.** Proportion of right choices on zero contrast trials as a function of the decoded Bayes-optimal prior, estimated from the neural activity: higher values of the decoded prior are associated with greater proportion of right choices. Ephys: decoding based on all neurons per probe, WFI: decoding based on all pixels. **f.** The corrected  $R^2$  for decoding the prior from neural activity are correlated with the corrected  $R^2$  for decoding the residual prior (Bayes-optimal prior minus

Bayes-optimal prior decoded from DLC features). This correlation implies that the Bayes-optimal prior decoded from neural activity cannot be explained simply by the motor features extracted by DLC.

### Post stimulus prior

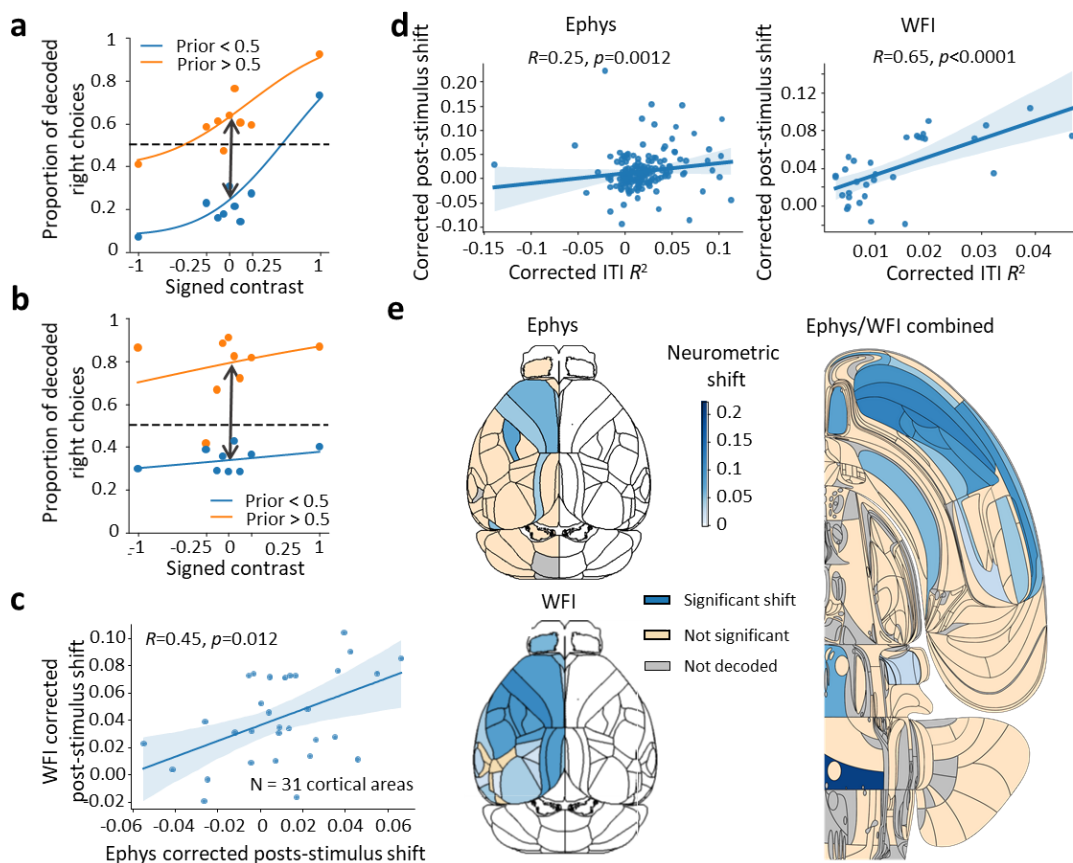
We also decoded the Bayes-optimal prior during the 100 ms interval after stimulus onset and found similarities in the encoding of the prior before and after stimulus onset. To avoid confounding the prior with the stimulus identity, two variables that are highly correlated (Spearman correlation  $R=0.45$ ,  $p<0.001$ ), we first trained a linear decoder of signed contrast from neural activity in each region. We used the output of this decoder to fit two neurometric curves (proportion of decoded right stimulus as a function of contrast, see Methods) conditioned on the Bayes-optimal prior being above or below 0.5. We then computed the vertical displacement of the fitted neurometric curves for zero contrast. If an area encodes the prior beyond the stimulus, we expect a shift between these two curves (see Fig. 3a for an example). Note that the same analysis can be performed during the ITI, though in this case the neurometric curves are expected to be flat (Fig. 3b) which is indeed what we observed (Fig. S9). This approach allows us to separate the encoding of the prior from the encoding of the stimulus, but it is possible that some of our results are related to the emergence of the animal's choices since the animals can respond in less than 100 ms on some trials (The International Brain Laboratory et al., 2021).

Using this approach, we found that we can detect the prior significantly from 20.7% (33/159) and 81.3% (26/32) of areas during the post-stimulus period for Ephys (Fig. S10a) and WFI respectively (using pseudosessions as described before for the null distribution). These percentages are smaller than for the ITI in part because this neurometric shift measure is less sensitive than direct decoding (in the ITI, only 10.6% of regions for Ephys and 93% for WFI are significant for the Bayes-optimal prior when using the neurometric shift on Ephys/WFI data, versus 29.5% and 100% respectively for straight decoding). As for the ITI prior decoding, we found that the Ephys and WFI post-stimulus shifts were correlated (Spearman correlation  $R=0.45$ ,  $p=0.012$ ,  $N=31$  Fig. 3c). Moreover, the post-stimulus neurometric shift is correlated with the  $R^2$  obtained in the same areas during the ITI period (Spearman correlation  $R=0.25$ ,  $p=1.2E-3$ ,  $N=159$  for Ephys,  $R=0.65$ ,  $p=4.9E-5$ ,  $N=32$  for WFI Fig. 3d). In other words, areas encoding the prior in the ITI also tend to do so during the post-stimulus period. This was confirmed by comparing the shifts during the post-stimulus and ITI periods, which were also found to be correlated (Fig. S10b).

We obtained similar results when merging the Ephys and WFI data into a single map (using Fisher's method to combine p-values across Ephys and WFI) and applying the Benjamini-Hochberg correction for multiple comparison (15.6% of significant regions, 25/160, Fig. 3e). Importantly, as observed during the ITI, areas encoding the prior were found at all levels of brain processing.

In addition, we explored whether regions encoding the stimulus also encoded the prior, as would be expected if these regions are involved in inferring the posterior distribution over the stimulus side. We found that the  $R^2$  for the stimulus decoding was indeed correlated with the corrected  $R^2$  for the Bayes-optimal prior decoding (Spearman correlation  $R=0.25$ ,  $p=0.006$ ,  $N=123$  regions from BWM analysis, (International Brain Lab et al., 2023) Fig. S11a). Moreover, among the 42 areas that were found to encode the stimulus significantly, 20 also encoded the prior significantly, including, once again, areas at all levels of brain processing (e.g., LGd, VISp, SCm, CP, MOs, MOp, ACAd, Fig. S11b).





**Figure 3:** Encoding of the prior across the brain during the post-stimulus period. **a.** Example of neurometric curves for post-stimulus period from the Caudate-Putamen. The double arrow shows the neurometric shift indicative of the prior encoding. **b.** Same as in **a** for the ITI period. **c.** Correlation between shifts from Ephys and WFI. Each dot corresponds to one cortical region. **d.** The post-stimulus shifts are correlated with the ITI  $R^2$  for both Ephys and WFI. **e.** Left: Comparison of corrected post-stimulus neurometric shift of Ephys and WFI for dorsal cortex. Right: Swanson map of corrected  $R^2$  averaged across Ephys and WFI for areas that have been deemed significant given both data sets (using Fisher's method for combining  $p$ -values), and after applying the Benjamini-Hochberg correction for multiple comparison.

### Decoding the action kernel prior

So far we have established that mice leveraged the block structure and that the Bayes-optimal prior can be decoded from the neural data at all levels of brain processing. However, it remains to be seen whether the animals truly compute the Bayes-optimal prior or, perhaps, use heuristics to compute a subjective, approximate, prior (Schaeffer et al., 2020).

To address this, we developed several behavioral models and used session-level Bayesian cross-validation followed by Bayesian model selection (Stephan et al., 2009) to identify the best fitting

model (see Methods). This analysis suggests that most animals on most sessions estimate what we will refer to as the action kernel prior, which is obtained by calculating an exponentially weighted average of recent past actions (Fig. 4a). The action kernel prior explains the mice's choices better than the Bayes-optimal prior and models of behavioral strategies that calculate an exponentially weighted average of recent stimuli (the 'stimulus kernel'), or assume a one step repetition bias, or the presence of positivity and confirmation biases (Palminteri & Lebreton, 2022) (see supp Fig. S12 and Supplementary Information). The decay constant of the exponential action kernel had a median of 5.6 trials across all animals (Fig. 4b, blue histogram), similar to the decay constant of recovery after block switches (5.16 trials, Fig. 1c). Remarkably, this is close to the value of the decay constant which maximizes the percentage of correct responses, given this (suboptimal) form of prior (Fig. 4b, orange curve). These curves are obtained by simulating the action kernel by varying the decay constant while keeping all other parameters at their best-fitting values.

If, as our behavioral analysis suggests, mice use the action kernel prior, then we should find that when we decode the prior inferred from the action kernel,  $R^2$  should be higher than when we decode the prior predicted by any other method. This is borne out by the data in both Ephys and WFI during the ITI (Fig. 4d). Unfortunately, however, and in contrast to the Bayes-optimal prior, we cannot determine which areas encode the action kernel prior significantly, because of the impossibility of generating a null distribution, as this would formally require having access to the exact statistical model of the animal behavior (see section 'Assessing the statistical significance of the decoding of the action kernel prior' in Methods).

To explore further whether neural activity better reflects the action kernel prior, as opposed to the stimulus kernel prior or the Bayes-optimal prior, we looked at changes in performance on zero contrast trials following correct and incorrect actions. When considering behavior within blocks, a decision-making agent using an action kernel prior should achieve a higher percentage of correct responses after a correct block-consistent action than an incorrect one, because, on incorrect trials, it updates the prior with an action corresponding to the incorrect stimulus side. Models simulating agents using either the Bayes-optimal prior or the stimulus kernel prior do not show this asymmetry since they perform their updates using the true stimulus, which can always be correctly inferred from the combination of action and reward (see also (Schaeffer et al., 2020)). Mice behavior showed the asymmetry in performance (Fig. 4c). Presumably, this asymmetry should also be present in the neural data. To test for this asymmetry, we decoded the Bayes-optimal prior from neural activity and simulated the animal's decision on each trial by selecting a choice according to whether the decoded prior is greater or smaller than 0.5 (i.e., assuming every trial has a zero contrast stimulus). We then asked whether the resulting sequence of hypothetical choices would show the asymmetry. If so, this is a property of the neural data since the predicted quantity, the Bayes-optimal prior, does not show the asymmetry. As shown in Fig. 4c, performance for both modalities, Ephys and WFI, were indeed higher following correct versus incorrect trials, thus strengthening our hypothesis that neural activity more closely reflects the action kernel prior.

Next, we tested the sensitivity of the decoded Bayes-optimal prior, estimated from neural activity, to previous actions (decoding the Bayes-optimal prior instead of the action kernel prior to allow us to test for statistical significance, see Methods). If the prior we estimate from neural activity reflects the subjective prior estimated from behavior, we should find that the neural prior is sensitive to the past 5-6

trials. Using an orthogonalization approach, we estimated the influence of past actions on the decoded Bayes-optimal prior and found that this influence extends at least to the past 5 trials in both Ephys and WFI (Fig. 4e, see Methods section titles ‘Orthogonalization’). A similar result was obtained when testing the influence of the past stimuli (Fig. 4e). These numbers are consistent with the decay constant estimated from behavior (5.6 trials). These results were obtained at the session level, by analyzing all available neurons. Furthermore, we analyzed single regions for which we had a large number of neurons recorded simultaneously (MOp, LP, GRN), or strong imaging signals (MOp, MOs, VISp). In all cases, we found that an asymmetry in the neural data following correct and incorrect choices as well as evidence that the Bayes-optimal prior decoded from these regions is influenced by the past 5-6 actions (Fig. S13).

These analyses also address one potential concern with our decoding approach. It is well known in the literature that animals keep track of the last action or last stimulus (Busse et al., 2011; Lak et al., 2020; Mendonça et al., 2020). It is therefore conceivable that our ability to decode the prior from neural activity is simply based on the encoding of the last action in neural circuits, which indeed provides an approximate estimate of the Bayes-optimal prior since actions are influenced by the prior (Fig. 1c). The fact that we observe an influence of the last 5-6 trials, and not just the last trial, rules out this possibility.

To test this even further, we estimated the temporal dependency of the WFI single pixel and Ephys single-unit activities on past actions directly and compared them to the behavioral sensitivity to past actions on the same sessions (both expressed in terms of neural learning rates, i.e., the inverse of the decay constants, see Methods). Note that this analysis tests whether the temporal dynamics of neural activity is similar to the temporal dynamics of the mouse behavior, defined by fitting the action kernel model, but without regressing first the neural activity against any prior. We found that the inverse decay constants of the neural activity are indeed correlated across sessions with the inverse decay constants obtained by fitting the action kernel model to behavior (Fig. 4f). Critically, this correlation goes away if we perform the same analysis using stimulus kernels instead of action kernels (Fig. S14). Moreover, these results established at a session level remained when accounting for the variability across mice (see Fig. S15).

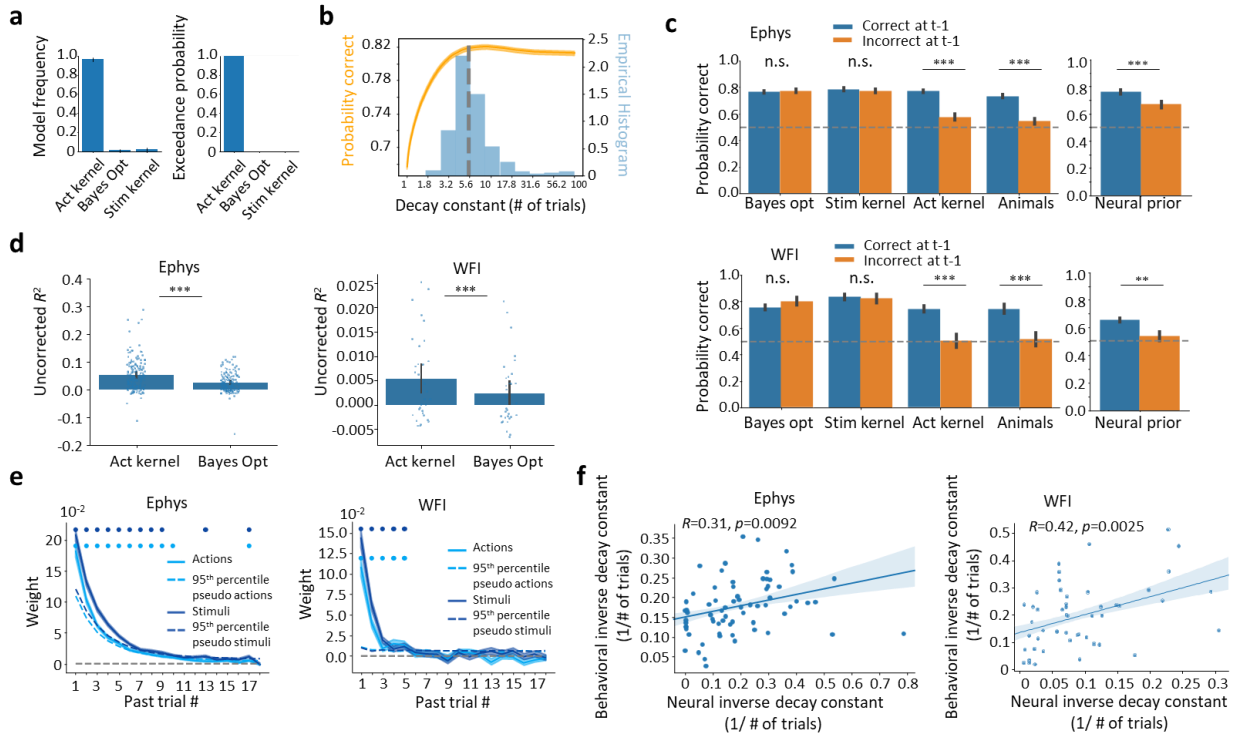


Figure 4: Action kernel prior. **a**. Model frequency (the posterior probability of the model given the subjects' data) and exceedance probability (the probability that a model is more likely than any other models (Stephan et al., 2009)) for three models of the mice's subjective prior. Across all sessions and all animals, the best model involves filtering the animal's last actions with an exponential kernel. Act kernel: action kernel model, Bayes Opt: Bayes-optimal model, Stim Kernel: stimulus kernel model. **b**. Blue: Histogram of the decay constant for the action kernel model across all animals. The median (dashed line) is 5.6 trials. Orange: probability of correct responses of the action kernel model as a function of the decay constant. The median decay time corresponds to the optimal performance. **c**. Performance on zero contrast trials conditioned on whether the previous action is correct or incorrect for various behavioral models and for the animals' behavior. Right: same analysis for a simulated agent using the Bayes-optimal prior decoded from neural data (Neural prior) to generate decisions. The decrease in performance between correct and incorrect previous trials for the neural prior suggests that the action kernel model best accounts for neural activity, which is consistent with behavior. Top: Ephys. Bottom: WFI. (\*\*\*)  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$  with a 2-tailed signed-rank paired Wilcoxon test) **d**. Uncorrected  $R^2$  is higher when decoding the action kernel prior during the ITI than when decoding the Bayes-optimal prior, for both Ephys and WFI (2 tailed-wilcoxon paired test for Ephys  $t=2148$ ,  $p=1.6E-13$ ,  $N=341$  sessions, and widefield  $t=81$ ,  $p=6.2E-4$ ,  $N=52$  sessions). **e**. Average weight of the previous actions (light blue) and previous stimuli (dark blue) on the predicted Bayes-optimal prior, estimated from neural activity. Ephys: results based on all neurons for a given Neuropixels probe, and averaged across probes/sessions. WFI: results based on all pixels and averaged across sessions. Filled circles indicate weights that are significantly different from zero. Dashed lines show the 95th percentile of the null distribution. **f**. Correlation between the neural inverse decay constants, obtained by estimating the temporal

dependency of the neural signals with respect to previous actions, and the behavioral inverse decay constants, obtained by fitting the action kernel model to the behavior. The neural and behavioral inverse decay constants are correlated for both Ephys (left, Pearson correlation  $R=0.31$ ,  $p=0.0092$ ,  $N=71$  sessions) and WFI (right, Pearson correlation  $R=0.42$ ,  $p=0.0025$ ,  $N=49$  sessions).

## Discussion

To summarize, our results indicate that the subjective prior is encoded, at least to some extent, at all levels of processing in the brain, including early sensory regions (e.g., LGd, VISp), associative regions (OFC, ACA, SCm) and motor regions (MOs, MOp, GRN). This finding brings further support to the hypothesis that information flows across the brain in a way that could support the sort of multidirectional inference apparent in Bayesian networks (Ackley et al., 1985; Berkes et al., 2011; Jardri et al., 2017; Lange & Haefner, 2022).

One might argue that what we call a ‘subjective prior’ might be better called ‘motor preparation’ in motor related areas, or a top-down ‘attentional signal’ in early sensory areas. Ultimately, though, what is important is not the term we use to refer to this signal, but rather that it has properties consistent with the subjective prior: 1- it is predictive of the animal’s choices, particularly on zero contrast trials (Fig. 2e), 2- it depends on previous actions (Fig. 4c), and 3- it reflects more than the last action or last stimulus (Fig. 4e), but depends instead on the past 5-6 choices (Fig. 4e). As we have seen, the signals we have recovered throughout the mouse brain fulfill all of these properties.

Our conclusions stand in partial contrast to the results of a companion paper (International Brain Lab et al., 2023) which reported the main portion of the Ephys data on which we also relied, and performed an initial set of consensus analyses of these data. That paper reported the encoding of the block prior along with several other variables such as stimulus, choice and feedback. In this other study, the block prior appeared to be the least prevalent variable, being found to be significantly encoded in only 6 regions when using a slightly different approach to the study of single cell responses (CENT3, CEA, CP, ORBV1 and SCm), and only in MOp when using a decoding approach. However, the analyses in the companion paper differ from the present one in several regards (see Methods section entitled ‘Difference with the Brain Wide Map prior results’ for a complete list of differences). In particular, we decoded the Bayes-optimal prior or the action kernel prior, while the other study focused on the true block prior, a quantity that is not available to the animal. Moreover, the present analysis considered more sessions, and a longer time interval, thus increasing our statistical power. And last but not least, we considered and combined the results of two recording modalities, while the companion paper relied exclusively on Ephys data.

There are several proposals in the literature as to how probability distributions might be encoded in neural activity. These include: linear probabilistic population codes (Ma et al., 2006), sampling based codes (Echeveste et al., 2020), other activity based codes (Dabney et al., 2020; Ganguli & Simoncelli, 2014; Sahani & Dayan, 2003; Schaeffer et al., 2020; Zemel et al., 1998), and the synaptic weights of neural circuits (Soltani & Wang, 2010). We note that our results are compatible with two requirements of linear probabilistic population codes (Ma et al., 2006; Walker et al., 2020): 1- the log odds of the Bayes-optimal prior is linearly decodable from neural activity (Fig. S5), and 2- changes in the Bayes-optimal prior from trial to trial are reflected in the population activity (Walker et al., 2020).

If the likelihood is also encoded with a linear probabilistic population code, having the prior in the same format would greatly simplify the computation of the posterior distribution, since it would simply require a linear combination of the neural code for the prior and likelihood. As it turns out, it is very likely that the likelihood indeed relies on a linear probabilistic population code. Indeed, the neural code for contrast, which is the variable that controls the uncertainty of the visual stimulus in our experiment, has been shown to be compatible with linear probabilistic population code (Ma et al., 2006).

Whether our results are also compatible with sampling based codes is more difficult to assess, as there is still a debate as to which aspects of neural activity correspond to a sample of a probability distribution (Echeveste et al., 2020; Haefner et al., 2016; Hoyer & Hyvärinen, 2003). Moreover, the fact that our prior follows a Bernoulli distribution, which is particularly simple, makes it harder to tease apart the various probabilistic coding schemes.

Ultimately, determining the exact nature of the neural code for the prior will require developing a neural model of Bayesian inference in a large, modular, loopy network - a pressing, remaining task. A critical foundation for this development is the remainder of the extensive data in the International Brain Laboratory brain-wide map (described in the companion paper, (International Brain Lab et al., 2023)). This provides a picture, at an unprecedented scale, of the neural processes underpinning decision-making, in which the prior plays such a critical part.

## Methods

The experimental methods concerning the International Brain Laboratory (IBL) data acquisition and preprocessing can be found in the methods of previous IBL publications. For a detailed account of the surgical methods for the headbar implants, see appendix 1 of (The International Brain Laboratory et al., 2021). For a detailed list of experimental materials and installation instructions, see appendix 1 of (Laboratory, International Brain, 2022). For a detailed protocol on animal training, see methods in (Laboratory, International Brain, 2022; The International Brain Laboratory et al., 2021). For details on the craniotomy surgery, see appendix 3 of (Laboratory, International Brain, 2022). For full details on the probe tracking and alignment procedure, see appendix 5 & 6 of (Laboratory, International Brain, 2022). The spike sorting pipeline used at IBL is described in detail in (Laboratory, International Brain, 2022). A detailed account of the widefield imaging data acquisition and preprocessing can be found in Christopher Krasniak's PhD thesis (Krasniak, 2022).

For the electrophysiological data, we used the 2022 IBL public data release (Laboratory, International Brain, 2023), which is extensively explored in IBL's latest preprint "Brain-wide neural activity during the IBL task". It is composed of 547 recordings from Neuropixels 1.0 probes. One or two probe insertions were realized over 354 sessions of the task, performed by a total of 115 mice. For the widefield calcium imaging data, we used another dataset consisting of 52 recordings of the dorsal cortex, realized over 52 sessions of the task, performed by a total of 6 mice. This second dataset is in the process of being made publicly available by the IBL team.

### Inclusion criteria for the analysis

**Criteria for trial inclusion.** All trials were included except when the animals did not respond to the stimulus (no movement, no response) or when the reaction time was shorter than 80 ms.

**Criteria for session inclusion.** All sessions were included except sessions with fewer than 400 trials (counting all trials, included and excluded) and sessions with fewer than 150 trials (counting only included trials).

**Criteria for neural recordings inclusion.** An insertion was included in the analysis if it had been resolved, that is, if histology clearly revealed the path of the probe throughout the brain, as defined in (Laboratory, International Brain, 2022). A neuron, identified during the spike sorting process, was included if it passed 3 quality control criteria (amplitude > 50  $\mu$ V ; noise cut-off < 20; refractory period violation). A region recorded along a probe was included in the analysis if there were at least 10 units that passed the quality control (QC). For widefield imaging, we use all the image pixels and have no criteria for pixel inclusion, nor region inclusion.

For the region-level analysis, after applying these criteria, we were left with 443 probe insertions (recorded over 308 sessions) for the Ephys dataset. For 83 insertions, none of the recorded regions passed the minimal units number criteria. For 5 insertions, the corresponding sessions did not pass the minimal number of trials criteria. Our region-level analysis spans 159 brain regions, defined by the Allen Common Coordinate Framework (Wang et al., 2020), recorded by at least one included insertion.

We define a region-insertion as the neurons recorded along a probe insertion, restricted to one region. Our analysis spans 931 region-insertions, which are aggregated across insertions to give results at the region level. For the embodiment analysis and the eye position analysis, we respectively looked at 242 region-insertions and 22 region-insertions, due to the additional criteria mentioned below.

For the region-level analysis where quality control criteria for neuron and region inclusion were relaxed meaning where the 3 previously quality control criteria as well as the minimum number of units per region constraint were relaxed (denoted also as QC=0 in the main text), we obtained a higher number of probe insertions (544), sessions (353) as well as brain regions included in the analysis (267).

For the session-level analysis, neurons along the whole probe were used and most of the probes recorded at least 10 units that passed the quality control. We obtained 519 QC probes insertions over 341 sessions.

All of the 52 widefield imaging sessions passed the minimal number of trials criteria and are thus included in the analysis. The imaging spans 32 regions of the dorsal cortex, which are part of the 159 regions decoded in the Ephys analysis, except for one area (AUDpo, auditory posterior area) .

**Criteria for the embodiment analysis.** Only sessions with available DLC features can be used for the embodiment prior analysis, which require access to body position and movement. For the Ephys data set, we analyzed the 106 sessions (out of 354) for which the DLC features met the quality criteria defined in (Laboratory, International Brain, 2022), and for which the other inclusion criteria were met. Wide field imaging sessions were excluded from this analysis as no video recordings were available.

**Criteria for the eye position analysis.** Reliable tracking of eye position from video recordings was not possible for some sessions because of video quality issues. Thus, we recovered reliable eye position signals from 21 out of the 27 of sessions in which we had recorded from either VISp or LGd, the two regions for which we analyzed the impact of eye position.

## Electrophysiological data

We used L1-regularized linear regression to decode the Bayes-optimal prior from the binned spike count data, with the scikit-learn function *sklearn.linear\_model.Lasso* (employing one regularization parameter  $\alpha$ ). We used L1 for Ephys because it is more robust to outliers, which are more likely to occur in single cell recordings, notably because of drift. The Bayes-optimal prior was inferred from the sequence of stimuli for each session (as described in Supplementary). Spike counts were obtained by summing the spikes across the decoding window for each included trial. If there were  $U$  units and  $T$  trials, this binning procedure resulted in a matrix of size  $U \times T$ . For the intertrial interval, the decoding window was (-600ms,-100ms) relative to the stimulus onset, and for the post-stimulus window, it was (0ms,+100ms) relative to stimulus onset. This decoding procedure yielded a continuous-valued vector of length  $T$ .

## Widefield imaging data

For the widefield calcium imaging data, we used L2-regularized regression as implemented by the scikit-learn function *sklearn.linear\_model.Ridge* (one regularization parameter  $\alpha$ ). We used L2 regularization instead of L1 for WFI data because L2 tends to be more robust to collinear features, which is the case across WFI pixels. We decoded the activity from the vector of the region's pixels for a specific frame of the data. The activity is the change in fluorescence intensity relative to the resting fluorescence intensity  $\Delta F/F$ . Data was acquired at 15 Hz. Frame 0 corresponds to the frame containing stimulus onset. For the intertrial interval, we use frame -2 relative to the stimulus onset, corresponding to an extremal time window of (-132ms,-66ms). For the post-stimulus interval, we use the frame +2, corresponding to an extremal time window of (+132ms,+198ms). If there are  $P$  pixels and  $T$  trials, this binning procedure results in a matrix of size  $P \times T$ .

## Reversal curves

To analyze mouse behavior around block reversals, we plot the reversal curves defined as proportion of correct choice as a function of trials, aligned to a block change (Fig. 1d). These are obtained by computing one reversal curve per mouse (pulling over sessions) and then averaging and computing the standard error across the mouse-level reversal curves. For comparison purposes, we also show the reversal curves for the Bayes-optimal model with a probability matching decision policy. We do not plot standard errors but standard deviations in this case as there is no variability across agents that we need to account for.

To formally assess differences between the mouse behavior and the agent that samples actions from the Bayes-optimal prior, we fit the following parametric function to the reversal curves:

$$p(\text{correct at trial } t) = (B + (A - B) \cdot e^{-t/\tau}) \cdot (t \geq 0) + B \cdot (t < 0)$$

with  $t = 0$  corresponding to the trial of the block reversal,  $\tau$  the decay constant,  $B$  the asymptotic performance and  $A$  the drop in performance right after a block change.

We fit this curve using the trials using only zero contrast trials, between the 5 pre-reversal trials and the 20 post-reversal trials. We restrict our analysis to the zero contrast trials to focus on trials where mice can only rely on block information to decide. This implies that we are only using a small fraction of the data. To be precise, across the 354 sessions, we have an average of 10 reversals per session, and the proportion of zero contrast trials is 11.1%. Fitting only on the zero contrast trials around reversals leads



us to use around 28 trials per session, which accounts for around 3% of the behavioral data (the average session length is 813 trials).

To make up for this limited amount of data, we use a jackknifing procedure for fitting the parameters. The procedure involves iteratively leaving out one mouse and fitting the parameters on the  $N-1=114$  zero contrast reversal curve of the held-in mice. Results of the jackknifing procedure are shown in Supplementary Fig. 1.

### Nested cross validation procedure

Decoding was performed using cross-validated, maximum likelihood regression. We used the scikit-learn python package to perform the regression (Pedregosa et al., 2011), and implemented a nested cross-validation procedure to fit the regularization coefficient.

The regularization parameter,  $\alpha$ , was determined via two nested five-fold cross-validation schemes (outer and inner). We first describe the procedure for the Ephys data. In the ‘outer’ cross-validation scheme, each fold is based on a training/validation set comprising 80% of the trials and a test set of the remaining 20% (random “interleaved” trial selection). The training/validation set is itself split into five sub-folds (‘inner’ cross-validation) using an interleaved 80-20% partition. Cross-validated regression is performed on this 80% training/validation set using a range of regularization weights, chosen for each type of dataset so that the bounds of the hyperparameter range are not reached (see table 1). The regularization weight selected with the inner cross-validation procedure on the training/validation set is then used to predict the target variable on the 20% of trials in the held-out test set. We repeat this procedure for each of the five ‘outer’ folds, each time holding out a different 20% of test trials such that, after the five repetitions, 100% of trials have a held-out decoding prediction. For widefield imaging, the procedure is very similar but we increased the number of outer folds to 50 and performed a leave-one-out procedure for the inner cross-validation. We did this because the number of features in widefield (number of pixels) is much larger than in Ephys (number of units): around 72 units in average in Ephys when decoding on a session-level from both probes after applying all quality criteria, vs around 2030 pixels on a session-level in widefield when decoding from the whole brain.

Furthermore, to average out the randomness in the outer randomization, we run this procedure 10 times, and take the mean of the obtained held-out predictions over the 10 runs. Each run uses a different random seed for selecting the interleaved train/validation/test splits. We take the average decoding score  $R^2$  across all runs to report the decoding score.

Type of decoding	regularization coefficient range $\alpha$
Ephys region-level	$\{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$
Ephys session-level	$\{10^{-2}, 10^{-1}, 1\}$
widefield region-level	$\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}\}$
widefield session-level	$\{10^{-2}, 10^{-1}, 1, 10, 100\}$

Ephys DLC features session-level	$\{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$
----------------------------------	---

**table 1** : hyperparameter grid for the cross-validation, different for each kind of decoding modality.

### Assessing statistical significance

Decoding a slow varying signal such as the Bayes-optimal prior from neural activity can easily lead to false positive results even when properly cross-validated. For instance, slow drift in the recordings can lead to spurious, yet significant, decoding of the prior if the drift is partially correlated with the block structure (Elber-Dorozko & Loewenstein, 2018; Harris, 2020). To control for this problem, we generated a null distribution of  $R^2$  values and determined significance with respect to that null distribution. This pseudosession method is described in detail in (Harris, 2020).

We will denote  $X \in R^{N \times p}$ , the aggregated neural activity for a session and  $Y \in R^N$  the Bayes-optimal prior (see ). Here,  $N$  is the number of trials and  $p$  the number of units. We generate the null distribution from "pseudosessions", i.e., sessions in which the true block and stimuli were resampled from the same generative process as the one used for the mice. This ensures that the time series of trials in each pseudosession shares the same summary statistics as the ones used in the experiment. For each true session, we generate  $M=200$  such pseudosessions, and use their resampled stimulus sequences to compute Bayes-optimal priors  $Y_i \in R^N$ , with  $i \in [1, M]$  the pseudosession number. We generate "pseudoscores"  $R_i^2 \in R$ ,  $i \in [1, M]$  by running the neural analysis on the pair  $(X, Y_i)$ . The neural activity  $X$  is independent of  $Y_i$  as the mouse didn't see  $Y_i$  but  $Y$ . Any predictive power from  $X$  to  $Y_i$  would arise from slow drift in  $X$  unrelated to the task itself. These pseudoscores  $R_i^2$  are compared to the actual score  $R^2$  obtained from the neural analysis on  $(X, Y)$  to assess statistical significance.

The actual  $R^2$  is deemed significant if it is higher than the 95th percentile of the pseudoscores  $R^2$  distribution. This test is used to reject the null hypothesis of no correlation between the Bayes optimal prior signal  $Y$  and the decoder prediction. We define the p-value of the decoding score as the quantile of  $R^2$  relative to the null distribution  $\{R_i^2, i \in [1, M]\}$ .

For each region of the brain that we recorded, we obtain a list of decoding p-values where a p-value corresponds to the decoding of the region's unit activity during one session. The number of recording sessions per region varies from 1 to 52 sessions (see Fig. S16). We use Fisher's method to combine the session-level p-values of a region into a single region-level p-value.

For effect sizes, we computed a corrected  $R^2$ , defined as the actual score  $R^2$  minus the median of the pseudoscores distribution,  $\{R_i^2, i \in [1, M]\}$ . The corrected  $R^2$  of a region is the mean of the corresponding sessions' corrected  $R^2$ .

### Controlling for region size when comparing decoding scores across Ephys and WFI

With WFI data, the activity signal of a region has always the same dimension across sessions, corresponding to the number of pixels. To control for the effect of region size on the region  $R^2$ , we performed linear regression across 31 recorded regions to predict the decoding  $R^2$  from the number of pixels per region. We found a significant correlation between  $R^2$  and the size of the regions (Fig S5a,  $R=0.88$   $p=4E-11$ ). To determine whether this accounts for the correlations between Ephys and WFI  $R^2$  correlation (Fig. 2d), we subtracted the  $R^2$  predicted by region's size from the WFI  $R^2$  and re-computed the correlation between Ephys  $R^2$  and these size-corrected WFI  $R^2$  (Fig. S5b).

### **Difference with the Brain Wide Map prior results**

The analyses presented in this paper differ from the prior decoding results presented in a companion paper (denoted BWM paper thereafter)(International Brain Lab et al., 2023) in the following respects:

1. We decoded the Bayes-optimal prior and action kernel prior, while the BWM paper considers the true block prior, a quantity that is not available to the animal given that block switches are not explicitly signaled to the animals.
2. The BWM paper only considered regions with two or more recording sessions, while we included all regions regardless of the number of recording sessions
3. Reaction times are assumed to be between 0.08 sec and 2 sec in the BWM paper while we included all trials with reaction time greater than 0.08 sec.
4. We decoded neural activity during the -600 to -100ms pre-stimulus onset, whereas the BWM paper focused on the interval -400 to -100ms.
5. We used linear regression instead of the logistic regression in the BWM paper.
6. We used 200 pseudosessions instead of 1000.
7. We treated data recorded simultaneously from two probes in the same region as two distinct independent datasets, while these were merged into a single dataset in the BWM analysis. To test the validity of this independence assumption, we correlated corrected  $R^2$  of identical regions across probes (same session, different probes, same region) and compared it with the correlations of corrected  $R^2$  of different regions across probes (same session, different probes, different regions). We have  $N=29$  session-regions that are present on both probes (same session, different probes, same region), leading to a Spearman correlation of  $R=0.34$ . The null hypothesis is that identical regions across probes are as independent as different regions across probes. To compute the null distribution, we randomly sampled  $K=10000$  sets of  $N=29$  corrected  $R^2$  of different regions across probes (same session, different probes, different regions), and computed the Spearman correlations of these  $K=10000$  sets. We then ask whether the  $R=0.34$  correlation is in the top 5% percentile of this null distribution. We find that  $R=0.34$  is in the top 32.4% percentile leading us to not reject the null hypothesis.

### **Proportion of right choices as a function of the decoded prior**

To establish a link between the decoded prior, estimated from the neural activity, and the mouse behavior, we plotted the proportion of right choices on zero contrast trials as a function of the decoded

Bayes-optimal prior. This analysis is performed at the session-level (in electrophysiology, concatenating all the units on the two probes when necessary). This means that, for each session for electrophysiology (or widefield), we decode the Bayes-optimal prior from all available neural activity on that session to define the decoded Bayes-optimal prior. We examine this on test trials, held-out during training, following the procedure described in the paragraph ‘Nested cross validation procedure’. For the main Fig. 2e, we then pool over decoded priors for all sessions, assign them to deciles and compute the associated proportion of right choices. In other words, we compute the average proportion of right choices on trials where the decoded prior is part of each decile.

To quantify the significance of this effect on a session-level (Supplementary Figure S7), we additionally performed a logistic regression predicting the choice (right or left) as a function of the decoded prior. Let  $j$  be the session number, we write the actions on that session  $a_t^j$  (with  $t$  the trial number) as a function of the decoded prior:

$$p(a_t^j = \text{right}) = 1 / \left[ 1 + \exp\left(-\left(\mu^j \cdot \widehat{Y}_t^j + c^j\right)\right) \right]$$

with  $\widehat{Y}_t^j$  the decoded Bayes-optimal prior,  $\mu^j$  the slope (coefficient of the logistic regression associated with the decoded prior) and  $c^j$  an intercept. The logistic regression fitting was performed with the default sklearn LogisticRegression function, which assumes a L2 regularization on weights with regularization strength  $C = 1$ .

To assess the statistical significance of these slopes,  $\mu^j$ , we generated null distributions of slopes over  $M$  pseudosessions (pseudosessions are defined in the paragraph above titled ‘Assessing statistical significance’). For each pseudosession, we computed the slope of the logistic regression between proportion of correct choices as a function of the decoded pseudo Bayes-optimal prior. The decoded pseudo Bayes-optimal prior was obtained by, first, computing the pseudo Bayes-optimal prior for each pseudosession, and then using the neural data from the original session to decode this pseudo Bayes-optimal prior. The percentage of correct choice was more complicated to obtain on pseudosession because it requires simulating the mice choices as accurately as possible. Since we don’t have a perfect model of the mice choices, we had to approximate this step with our best model, i.e., the action kernel model. We used the action kernel model fitted to the original behavior session and simulated it on each pseudosession to obtain the actions on each trial of the pseudosessions.

From the set of decoded pseudo Bayes-optimal priors and pseudoactions, we obtained  $M$  pseudoslopes  $\mu_i^j$ ,  $i = 1 \dots M$  using the procedure described above. Because the mouse did not experience the pseudosessions or perform the pseudoactions, any positive coefficient  $\mu_i^j$  has to be the result of spurious correlations. Formally, to assess significance, we ask if the mean slope  $\left( \mu = 1/J \cdot \sum_j \mu^j; j \in [1, J] \right)$  is within the 5% top percentile of the averaged pseudoslopes:  $\left\{ \mu_i; \mu_i = 1/J \cdot \sum_j \mu_i^j; i \in [1, M] \right\}$ . Fig. S7 shows this set of  $M$  averaged pseudoslopes as a histogram. The red vertical dashed line is the average slope  $\mu$ .

When applying this null-distribution procedure in Ephys and WFI data, we find that the pseudoslopes in Ephys are much more positive than in WFI. This is due to the fact that spurious correlations in Ephys are induced by drift in the Neuropixels probes, while WFI data barely exhibits any drift.

## Neurometric Curves

We use the same decoding pipeline described for the Bayes-optimal prior decoding to train a linear decoder of the signed contrast from neural activity in each region, for the intertrial interval  $[-600, -100ms]$  and post-stimulus  $[0, 100ms]$  intervals. There are 9 different signed contrasts  $\Gamma = \{-1, -0.25, -0.125, -0.0625, 0, 0.0625, 0.125, 0.25, 1\}$  where the left contrasts are negative and the right contrasts are positive. Given a session of  $N$  trials, we denote  $c \in R^N$  the sequence of signed contrasts and  $\hat{c} \in R^N$  the decoder output given the neural activity  $X$ . For every contrast  $c$ , we compute the proportion of decoded right stimuli conditioned on the Bayes-optimal prior being higher or lower than 0.5. Given the trial indexes sets of a session  $I_c^{low}$  and  $I_c^{high}$  corresponding to the trials with signed contrast  $c$  and a Bayes-optimal prior lower or higher than 0.5 respectively, we compute the proportions  $P_c^{low} = \#\{\hat{c}_i > 0; i \in I_c^{low}\} / \#I_c^{low}$  and  $P_c^{high} = \#\{\hat{c}_i > 0; i \in I_c^{high}\} / \#I_c^{high}$ . We fit a *low prior* curve to  $\{(c, P_c^{low})\}_{c \in \Gamma}$  and a *high prior* curve to  $\{(c, P_c^{high})\}_{c \in \Gamma}$ , which we call neurometric curves. We use an *erf()* function from 0 to 1 with two lapse rates for the curves fit to obtain the neurometric curve:

$$f(c) = \gamma + (1 - \gamma - \lambda) * (erf((c - \mu)/\sigma) + 1)/2$$

Where  $\gamma$  is the low lapse rate,  $\lambda$  is the high lapse rate,  $\mu$  is the bias (threshold) and  $\sigma$  is the rate of change of performance (slope). We use the psychofit toolbox to fit the neurometric curves using maximal likelihood estimation (<https://github.com/cortex-lab/psychofit>). Finally, we estimate the vertical displacement of the fitted neurometric curves for the zero contrast  $f^{high}(c = 0) - f^{low}(c = 0)$ , which we refer to as the neurometric shift.

We use the pseudosession method to assess the significance of the neurometric shift, by constructing a neurometric shift null distribution.  $M=200$  pseudosessions are generated with their signed contrast sequences, which are used as target to linear decoder on the true neural activity. We fit neurometric curves to the pseudosessions decoder outputs, conditioned on the Bayes-optimal prior inferred from the pseudosessions contrast sequences.

## Embodiment

Video data of two cameras were used to extract 7 behavioral variables which could potentially be modulated according to the mice's subjective prior (see (Laboratory, International Brain, 2022): licking, whisking left & right, wheeling, nose position and paw position left & right. If we are able to significantly decode the Bayes-optimal prior from these behavioral variables during the  $[-600 ms, -100 ms]$  inter-trial interval, we say that the subject displays an embodiment of the prior. For the decoding, we used L1-regularized maximum likelihood regression with the same cross-validation scheme used for neural data (see paragraph Nested Cross-Validation procedure). Sessions and trials are subjected to the

same QC as for the neural data, so that we decode the same sessions and the same trials as the Ephys session-level decoding. For the session of  $T$  trials, the decoder input is a matrix of size  $T \times 7$  and the target is the inferred Bayes-optimal prior. Every behavioral variable is binned into 0.02s bins, which we average over the ITI. We use the pseudosession method to assess the significance of the DLC features decoding score  $R^2$ .

DLC features	Description
paws position	Euclidean distance of the DLC-tracked paws to a camera frame corner.
nose position	x-position of the DLC-tracked nose using the left camera.
wheeling	defined by interpolating the wheel position at 5Hz and computing the motion magnitude.
licking	Left & right edge of the tongue are DLC-tracked using both lateral cameras. A lick is defined to have occurred in a frame if the difference for either coordinate to the subsequent frame is larger than 0.25 times the standard deviation of the difference of this coordinate across the whole session. The timestamps of the detected licks were binned into 0.02 sec bins, resulting in a lick activity metric.
whisking	defined by the motion energy of the whisker pad area filmed by camera (i.e. the mean across pixels of the absolute value of the difference of gray scales between adjacent frames); respectively by using the left & right camera (150Hz).

To investigate the embodiment of the Bayes-optimal prior signal, we compare session-level decoding of the prior signal from DLC regressors to insertion-level decoding of the prior signal from the neural activity of each region along the insertion. For the 106 sessions which pass the DLC quality criteria, we looked at 275 individual decoding from the regions along the corresponding 154 insertions.

### **DLC Residual analysis**

The DLC prior residual signal is the part of the prior signal which is not explained away by the DLC decoding, defined as the prior signal minus the prediction of the DLC decoding. We decode this DLC prior residual signal from the neural activity, using the same linear decoding schemes as previously described.

### **Eye position decoding**

Video data from the left camera were used to extract the eye position variable, a 2D-signal corresponding to the position of the center of the mouse pupil relative to the video border. The camera as well as the mouse's head are fixed. DeepLabCut was not able to achieve sufficiently reliable tracking of the pupils; therefore we used an improved pose estimation algorithm (Biderman et al., 2023), trained on the same labeled dataset used to train DeepLabCut. For the decoding, we used L2-regularized maximum likelihood regression with the same cross-validation scheme used for neural data, during the [-600 ms, -100 ms] inter-trial interval.

The eye-position prior residual signal is the part of the prior signal which is not explained away by the eye position decoding, defined as the prior signal minus the prediction of the eye position decoding. We decode this eye position prior residual signal from the neural activity of early visual areas (LGd & VISp), using the same linear decoding schemes as previously described.

## Behavioral models

To determine the behavioral strategies used by the mice, we developed several behavioral models and used Bayesian model comparison to identify the one that fits best. We considered three types of behavioral models which differ as to how the integration across trials is performed (how a prior probability, that the stimulus will be on the right side, is estimated based on history). Within a trial, all models compute the posterior distribution by taking the product of a prior and a likelihood function (probability of the noisy contrast given the stimulus side, see Supplementary information).

Among the three types of models of the prior, the first, called the *Bayes-optimal model*, assumes knowledge of the generative process of the blocks. Block lengths are sampled as follows:

$$p(l_k = N) \propto \exp(-N/\tau) \cdot \mathbb{1}[20 \leq N \leq 100]$$

With  $l_k$  the length of block  $k$  and  $\mathbb{1}$  the indicator function. Block lengths are thus sampled from an exponential distribution with parameter  $\tau = 60$  and constrained to be between 20 and 100 trials. When block  $k - 1$  comes to an end, the next block  $b_k$ , with length  $l_k$ , is defined as a “right” block (where stimulus will appear more frequently on the right) if block  $b_{k-1}$  was a “left” block (where stimulus appeared more frequently on the left) and conversely. During “left” blocks, stimulus will be on the left side with probability  $\gamma = 0.8$  (and similarly for “right” blocks). With  $s_t$  the stimulus side at trial  $t$ , the prior probability the stimulus will appear on the right at trial  $t$ ,  $p(s_t | s_{1:(t-1)})$  is obtained through a likelihood recursion (Scott, 2002).

The second, called the *stimulus kernel model* (Norton et al., 2019), assumes that the prior is estimated by integrating previous stimuli with an exponentially decaying kernel. With  $s_{t-1}$  the stimulus side at trial  $t - 1$ , the prior probability that the stimulus will appear on the right  $\pi_t$  is updated as follows:

$$\pi_t = (1 - \alpha) \cdot \pi_{t-1} + \alpha \cdot \mathbb{1}[s_{t-1} = \text{right}]$$

with  $\pi_{t-1}$  the prior at trial  $t - 1$  and  $\alpha$  the learning rate. The learning rate governs the speed of integration: the closer  $\alpha$  is to 1, the more weight is given to recent stimuli  $s_{t-1}$ .

The third, called the *action kernel model*, is similar to the *stimulus kernel model* but assumes an integration over previous chosen actions with, again, an exponentially decaying kernel. With  $a_{t-1}$  the action at trial  $t - 1$ , the prior probability that the stimulus will appear on the right  $\pi_t$  is updated as follows:

$$\pi_t = (1 - \alpha) \cdot \pi_{t-1} + \alpha \cdot \mathbb{1}[a_{t-1} = \text{right}]$$

For the *Bayes-optimal* and *stimulus kernel* models, we additionally assume the possibility of capturing a simple autocorrelation between choices with an immediate repetition bias or choice- and outcome-dependent learning rate (Palminteri & Lebreton, 2022; Sugawara & Katahira, 2021). See Supplementary Information for more details on model derivations.

## Model comparison

To perform model comparison, we implement a session-level Bayesian cross validation procedure, where, for each mouse, we hold out one session  $i$  and fit the model on the held-in sessions. For each mouse, given a held-out session  $i$ , we fit each model  $k$  to the actions of held-in sessions, denoted here as  $A^{\setminus i}$  and obtain the posterior probability,  $p(\theta_k | A^{\setminus i}, m_k)$ , over the fitted parameters  $\theta_k$  through an adaptive Metropolis-Hastings (M-H) procedure (Andrieu & Thoms, 2008).  $\theta_k$  will typically include sensory noise parameters, lapse rates and the learning rate (for *stimulus* and *action Kernel* models) - see Supplementary Information for the formal definitions of these parameters. Let  $\{\theta_k^n; n \in [1, N_{MH}]\}$  be the  $N_{MH}$  samples obtained with the M-H procedure for model  $k$  (after discarding the burn-in period). We then compute the marginal likelihood of the actions on the held-out session, denoted here as  $A^i$ .

$$p(A^i | A^{\setminus i}, m_k) = \int p(A^i, \theta_k | A^{\setminus i}, m_k) \theta_k = \int p(A^i | \theta_k, m_k) p(\theta_k | A^{\setminus i}, m_k) d\theta_k \approx \frac{1}{N_{MH}} \sum_n p(A^i | \theta_k^n, m_k)$$

For each subject, we obtain a score per model  $k$  by summing over the log-marginal likelihoods  $\log p(A^i | A^{\setminus i}, m_k)$ , holding out one session at a time. Given these subject-level log-marginal likelihood scores, we perform Bayesian model selection (Stephan et al., 2009) and report the model frequencies (the expected frequency of the  $k$ -th model in the population) and the exceedance probabilities (the probability that a particular model  $k$  is more frequent in the population than any other considered model).

## Assessing the statistical significance of the decoding of the action kernel prior

Given that the action kernel model better accounts for the mice's behavior, it would be desirable to assess the statistical significance of the decoding of the action kernel prior. Crucially, since assessing significance involves a null hypothesis (the neural activity is independent of the prior), a careful and rigorous construction of the corresponding null distribution is key.

For the Bayes optimal prior decoding, constructing the null distribution is straightforward. It requires that we generate stimulus sequences with the exact same statistics as those experienced by the mice. We do this by simulating the same generative process used to generate the stimulus during the experiment, yielding what we called pseudosessions in previous sections.



However, for the action kernel prior (and contrary to the Bayes optimal model), we also need to generate action sequences with the same statistics as those generated by the animals. In turn, this would require a perfect model of how the animals make decisions. Since we lack such a model, we would need to come up with an approximation. There are multiple approximations that we could use, including:

- 1- Synthetic sessions, in which we use the action kernel model, using the parameters fitted to each mouse on each session, to generate fake responses. However, the action kernel model is not a perfect model of the animal's behavior, it is merely the best model we have among the ones we have tested. Additionally, there could be some concerns about the statistical validity of using a null distribution, which assumes that the action kernel is the perfect model when testing for the presence of this same model in the mouse's neural activity
- 2- Imposter sessions, in which we use responses from other mice. However, other animals are most unlikely to have used the exact same model/parameters as the mouse we are considering. This implies that the actions in these imposter sessions do not have the same statistics as the decoded session. There is indeed a large degree of between-session variability, as can be seen from the substantial dispersion in the fitted action kernel decay constants shown in Fig. 4b.
- 3- Shifted sessions, in which we decode the action kernel prior on trial  $M$ , using the recording on trial  $M+N$ , with periodic boundaries for the 'edges'. The problems here are two-fold. First,  $N$  must be chosen large enough such that the block structure of the shifted session is independent of the block structure of the non-shifted session. Because blocks are about 50 trials long,  $N$  must be large for the independence assumption to hold. This adds a constraint on the number of different shifted sessions that we can generate, leading to a poor null distribution with little diversity (made from only a few different shifted sessions). Second, it has been shown that there is within-session variability (Ashwood et al., 2022) such that when  $N$  is chosen large, we can not consider the shifted actions to have the same statistics as the non-shifted actions.

There may be other options. However, since they would all rely on approximations, the degree of statistical inaccuracy associated with their use is unclear. We would not even know which one to favor, as it is hard to establish the quality of the approximations. Overall, we have access to the exact generative process to construct the null distribution for the Bayes optimal prior, versus only approximations for the action kernel prior. As a result, we decided to err on the side of caution and report statistical significance for the Bayes optimal prior decoding.

### **Orthogonalization**

To assess the dependency on past trials of the decoded Bayes-optimal prior from neural activity, we performed stepwise linear regression as a function of the previous actions (or previous stimuli). The Bayes-optimal prior was decoded from neural activity on a session-level, thus considering the activity from all accessible cortical regions in WFI and all units in Ephys (concatenating over probe insertions).

The stepwise linear regression involved the following steps. We started by linearly predicting the decoded Bayes-optimal prior on trial  $t$  from the previous action (action on trial  $t-1$ ), which allows us to compute a first-order residual, defined as the difference between the decoded neural prior and the decoded prior predicted by the last action. We then used the second-to-last action (action at trial  $t-2$ ) to predict the first-order residual, to then compute a second-order residual. Next, we predicted the second-order residual with the third-to-last action and so on. We use this iterative stepwise procedure in order to take into account possible autocorrelations in actions.

The statistical significance of the regression coefficients is assessed as follows. Let us denote  $N_j$  the number of trials of session  $j$ ,  $Y^j \in \mathbb{R}^{N_j}$  the decoded Bayes-optimal prior, and  $X^j \in \mathbb{R}^{N_j \times K}$  the chosen actions, where  $K$  is the number of past trials considered in the stepwise regression. When running the stepwise linear regression, we obtain a set of weights  $\{W_k^j, \text{ with } k \in [1, K]\}$ , with  $W_k^j$  the weight associated with the  $k^{\text{th}}$ -to-last chosen action. We test for the significance of the weights for each step  $k$ , using as a null hypothesis that the weights associated with the  $k^{\text{th}}$ -to-last chosen action are not different from weights predicted by the null distribution.

To obtain a null distribution, we followed the same approach as in the section entitled *Proportion of right choices as a function of the decoded prior*. Thus, we generated decoded pseudo Bayes-optimal priors and pseudoactions. For each session, these pseudovariables are generated in the following way: first, we fitted the action kernel model (our best fitting-model) to the behavior of session  $j$ . Second, we generated  $M$  pseudosessions (see the *Assessing the statistical significance* section). Lastly, we simulated the fitted model on the pseudosessions to obtain pseudoactions. Regarding the decoded pseudo Bayes-optimal priors, we first infer with the Bayes-optimal agent, the Bayes-optimal prior of the pseudosessions, and second, we decoded this pseudoprior with the neural activity. For each session  $j$  and pseudo  $i$ , we have generated a decoded pseudo Bayes-optimal prior  $Y_i^j$  as well as pseudoactions  $X_i^j$ . When applying the stepwise linear regression procedure to the couple  $(X_i^j, Y_i^j)$ , we obtain a set of pseudoweights  $\{W_{k,i}^j, \text{ with } k \in [1, K]\}$ . Because the mouse didn't experience the pseudosessions or perform the pseudoactions, any non-zero coefficients  $W_{k,i}^j$  must be the consequence of spurious correlations. Formally, to assess significance, we ask if the average of the coefficients over session

$$W_k = 1/N_{\text{sessions}} \cdot \sum_{j=1}^{N_{\text{sessions}}} W_k^j \quad \text{is within the 5\% top percentile of}$$

$$\left\{ W_{k,i}^j; W_{k,i} = 1/N_{\text{sessions}} \cdot \sum_{j=1}^{N_{\text{sessions}}} W_{k,i}^j; i \in [1, M] \right\}.$$

The statistical significance procedure when predicting the decoded Bayes-optimal prior from the previous stimuli is very similar to the one just described for the previous actions. The sole difference is that, for this second case, we do not need to fit any behavioral model to generate pseudostimuli. Pseudostimuli for sessions  $j$  are defined when generating the  $M$  pseudosessions. Pseudoweights are then obtained by running the stepwise linear regression predicting the decoded pseudo Bayes-optimal prior from the pseudostimuli. Formal statistical significance is established in the same way as for the previous actions case.

When applying this null-distribution procedure to Ephys and WFI, we find that the strength of spurious correlations (as quantified by the amplitude of pseudoweights  $W_{k,i}$ ) for Ephys is much greater than for WFI data. This is due to the fact that spurious correlations in electrophysiology are mainly produced by drift in the Neuropixels probes, which is minimized in WFI.

### **Behavioral signatures of the action kernel model**

To study why the Bayesian model selection procedure favors the action kernel model, we sought behavioral signatures that can be explained by this model but not the others. Since the action kernel model integrates over previous actions (and not stimuli sides), it is a self-confirmatory strategy. This means that if an action kernel agent was incorrect on a block-conformant trial (trials where the stimulus is on the side predicted by the block prior), then it should be more likely to be incorrect on the subsequent trial (if it is also block-conformant). Other models integrating over stimuli, such as the Bayes-optimal or the stimulus Kernel model are not more likely to be incorrect following an incorrect trial, because they can use the occurrence or non-occurrence of the reward to determine the true stimulus side, which could then be used to update the prior estimate correctly. To test this, we analyzed the proportion correct of each session at trial  $t$ , conditioned on whether it was correct or incorrect at trial  $t - 1$ . To isolate the impact of the last trial, and not previous trials or other factors such as block switches and structure, we restricted ourselves to:

- zero contrast trials
- on trials that are least 10 trials from the last reversal
- trial  $t$ ,  $t - 1$ , and  $t - 2$ , had stimuli which were on the “expected”, meaning “block-conformant” side
- on trial  $t - 2$ , the mouse was correct, meaning that it chose the block-conformant action

### **Neural signature of the action kernel model from the decoded Bayes-optimal prior**

To test if the behavioral signature of the action kernel model discussed in the previous section is also present in the neural activity, we simulated an agent whose decisions are based on the decoded Bayes-optimal prior and tested whether this agent also shows the same action kernel signature. The decoded Bayes-optimal prior was obtained by decoding the Bayes-optimal prior from the neural activity (see *Nested cross validation procedure* section) on a session-level basis, considering all available widefield pixels or electrophysiology units and concatenating units across probe insertions when necessary.

Note that if the decoded Bayes-optimal agent exhibits the action kernel behavioral signature, this must be a property of the neural activity since the Bayes-optimal prior on its own cannot produce this behavior.

The agent is simulated as follows. Let us denote  $Y \in R^N$  the Bayes-optimal prior with  $N$  is the number of trials. When performing neural decoding of the Bayes-optimal prior  $Y$ , we obtain a decoded Bayes-optimal prior  $\hat{Y}$ . We define an agent which, on each trial, greedily selects the action predicted by

the decoded Bayes-optimal prior  $\hat{Y}$ , meaning that the agent chooses right if  $\hat{Y} > 0.5$ , and left otherwise.

On sessions which significantly decoded the Bayes optimal prior, we then test whether the proportion of correct choices depends on whether the previous trial was correct or incorrect. We do so at the sessions level, applying all but one criteria of the behavioral analysis described previously in the *Behavioral signatures of the action kernel model* paragraph:

- on trials that are least 10 trials from the last reversal
- trial  $t$ ,  $t - 1$ , and  $t - 2$ , had stimuli which were on the “expected”, meaning on the “block-conformant” side
- on trial  $t - 2$ , the mouse was correct, meaning that it chose the block-conformant action

Note that, given that the neural agent uses the pre-stimulus activity to make its choice, we do not need to restrict ourselves to zero contrast trials.

## Neural Decay Rate

To estimate the temporal dependency of the neural activity in Ephys and WFI, we assume that the neural activity is the result of an action kernel (or stimulus kernel) integration and fit the learning rate (inverse decay rate) of the kernel to maximise the likelihood of observing the neural data.

We first describe the fitting procedure for widefield imaging data. Given a session, let us call  $X_n^t$  the widefield calcium imaging activity of the  $n$ -th pixel for trial  $t$ . Similarly to the procedure we used for decoding the Bayes-optimal prior, we take the activity at the second-to-last frame before stimulus onset. We assume that  $X_n^t$  is a realization of Gaussian distribution with mean  $Q_n^t$  and with standard deviation  $\sigma_n$ ,  $X_n^t \sim N(Q_n^t, \sigma_n)$ .  $Q_n^t$  is obtained through an action kernel (or stimulus kernel) integration process:

$$Q_n^t = (1 - \alpha_n) \cdot Q_n^{t-1} + \alpha_n \cdot \zeta_n \cdot a_{t-1}$$

with  $\alpha_n$  the learning rate,  $a_{t-1} \in \{-1, 1\}$  the action at trial  $t - 1$  and  $\zeta_n$  a scaling factor.  $\alpha_n$ ,  $\zeta_n$  and  $\sigma_n$  are found by maximizing the probability of observing the widefield activity  $p(X_n^{1:T} | a^{1:T}; \alpha_n, \zeta_n, \sigma_n)$ , with  $1:T = \{1, 2, \dots, T\}$  and  $T$  the number of trials in that session.

For the electrophysiology now, let us call  $X_n^t$  the neural activity of unit  $n$  at trial  $t$ . Similarly to what we did when decoding the Bayes-optimal prior, we take the sum of the spikes between -600 and -100ms from stimulus onset. We assume here that  $X_n^t$  is a realization of a Poisson distribution with parameter  $Q_n^t$ ,  $X_n^t \sim \text{Poisson}(Q_n^t)$ .  $Q_n^t$  is obtained through an action kernel (or stimulus kernel) integration process:

$$Q_n^t = (1 - \alpha_n) \cdot Q_n^{t-1} + \alpha_n \cdot \zeta_n^{a_{t-1}}$$

with  $\alpha_n$  the learning rate and  $\zeta_n^{a^{t-1}}$  scaling factors, one for each possible previous action.  $\alpha_n$ ,  $\zeta_n^1$  and  $\zeta_n^{-1}$  are found by maximizing the probability of observing the Ephys activity  $p(X_n^{1:T} | a^{1:T}; \alpha_n, \zeta_n^1, \zeta_n^{-1})$ , with  $1:T = \{1, 2, \dots, T\}$  and  $T$  the number of trials in that session. For electrophysiology, we add constraints on the units we will consider, specifically, we only consider units 1- whose median (pre-stim summed) spikes is not 0, 2- there must be at least 1 spike every 5 trials and 3- the distribution of (pre-stim summed) spikes must be different when the Bayes-optimal prior is greater versus lower than 0.5 (significance is asserted when the p-value of a Kolmogorov Smirnov test is below 0.05).

To only consider units (or pixels) which are likely to reflect the subjective prior, we restrict our analysis to units (or pixels) which are part of regions-insertions that significantly decode the Bayes-optimal prior (16% of the total number of region-insertions in Ephys and 61% in widefield). Then, to obtain a session-level neural learning rate, we average across pixel-level or unit-level learning rates. To compare neural and behavioral temporal timescales, we then correlate the session-level neural learning rate with the behavioral learning rate, obtained by fitting the action kernel to the behavior.

In both Ephys and WFI, when considering that the neural activity is a result of the stimulus kernel, the calculations are all identical except that one must replace actions  $a^{1:T}$  by stimuli side  $s^{1:T}$ .

This analysis (presented in Fig. 4f) makes the assumption that sessions can be considered as independent from another - assumption, which can be questionable given that we have 354 sessions across 115 mice in electrophysiology and 52 sessions across 6 mice in widefield. To test the presence of the correlation between neural and behavioral timescales while relaxing this assumption, we developed a hierarchical model which takes into account the two types of variability, within mice and within sessions given a mouse. This model defines session-level parameters, which are sampled from mouse-level distributions, which are themselves dependent on population-level distributions. See Supplementary information for the exact definition of the hierarchical model. This hierarchical approach confirmed the session-level correlation between neural and behavioral timescales (see supplementary Fig. S15).

## Data availability

All the data that support the findings of the present study are available at <https://int-brain-lab.github.io/iblenv>. Users are allowed to distribute, remix, adapt, and build upon the material in any medium or format, so long as attribution is given to the creator (data license CC-BY).

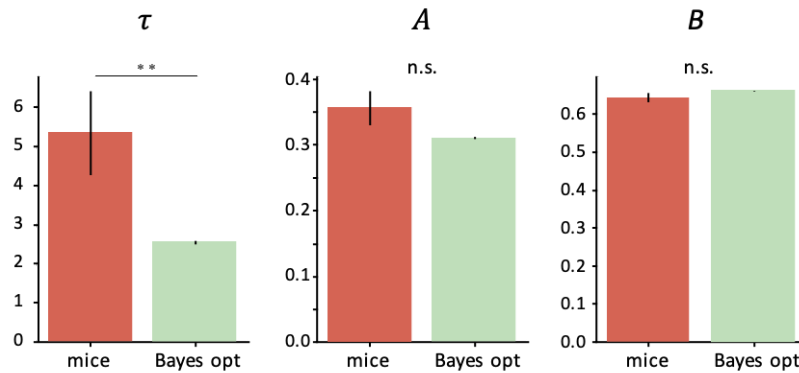
## Acknowledgements

This work was supported by grants from the Wellcome Trust (216324), the Simons Foundation, The National Institutes of Health (NIH U19NS12371601), the National Science Foundation (NSF 1707398), the Gatsby Charitable Foundation (GAT3708), and by the Max Planck Society and the Humboldt Foundation. We would like to thank the University of Geneva for providing computational resources and support that contributed to these research results.

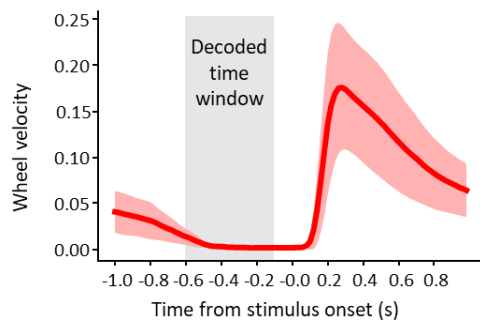
## Competing interests

The authors declare no competing interests

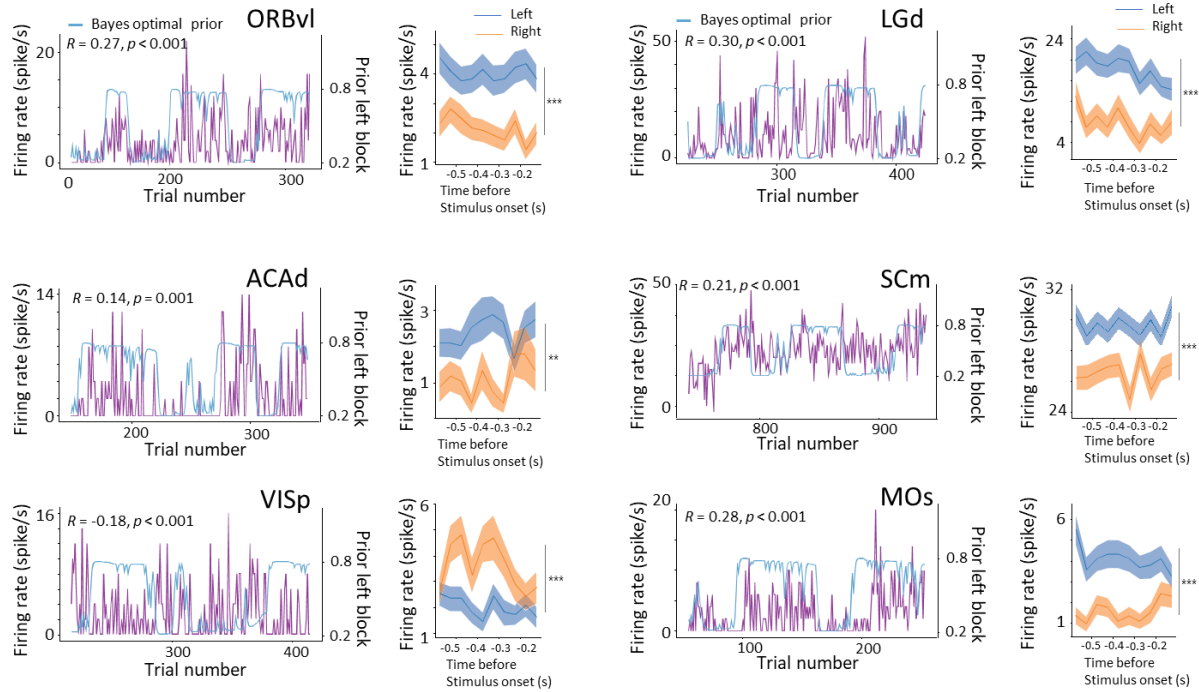
## Supplementary Figures



**Figure S1.** Histograms of the decay constant ( $\tau$ ), amplitude ( $A$ ) and asymptote ( $B$ ) of the zero contrast reversal curves across all mice. The parameters are obtained by fitting the following parametric curve  $p(\text{correct at trial } t) = B$  on the zero contrast pre-reversal trials (the 5 trials before a block switch) and  $p(\text{correct at trial } t) = B + (A - B) \cdot e^{-t/\tau}$  on the zero contrast post-reversal trials (the 20 trials after a block switch).  $\tau$  reflects the reversal timescale. To make up for the limited amount of available zero contrast reversal trials, we fit these curves using a jackknife procedure (see Methods). Bars and error bars indicate jackknife means  $\pm$  SEM (jackknifing was applied on  $N = 115$  mice). Mice have a significantly longer mean recovery decay constant than the Bayes-optimal observer (5.16 vs 2.46,  $t_{114}=2.94$ ,  $p=0.004$ ), while the other parameters are not significantly different. (for  $A$ :  $t_{114}=1.7$ ,  $p=0.09$  and for  $B$ :  $t_{114}=-0.49$ ,  $p=0.63$ ) (\*  $p<0.05$ , \*\*  $p<0.01$ , n.s. not significant)

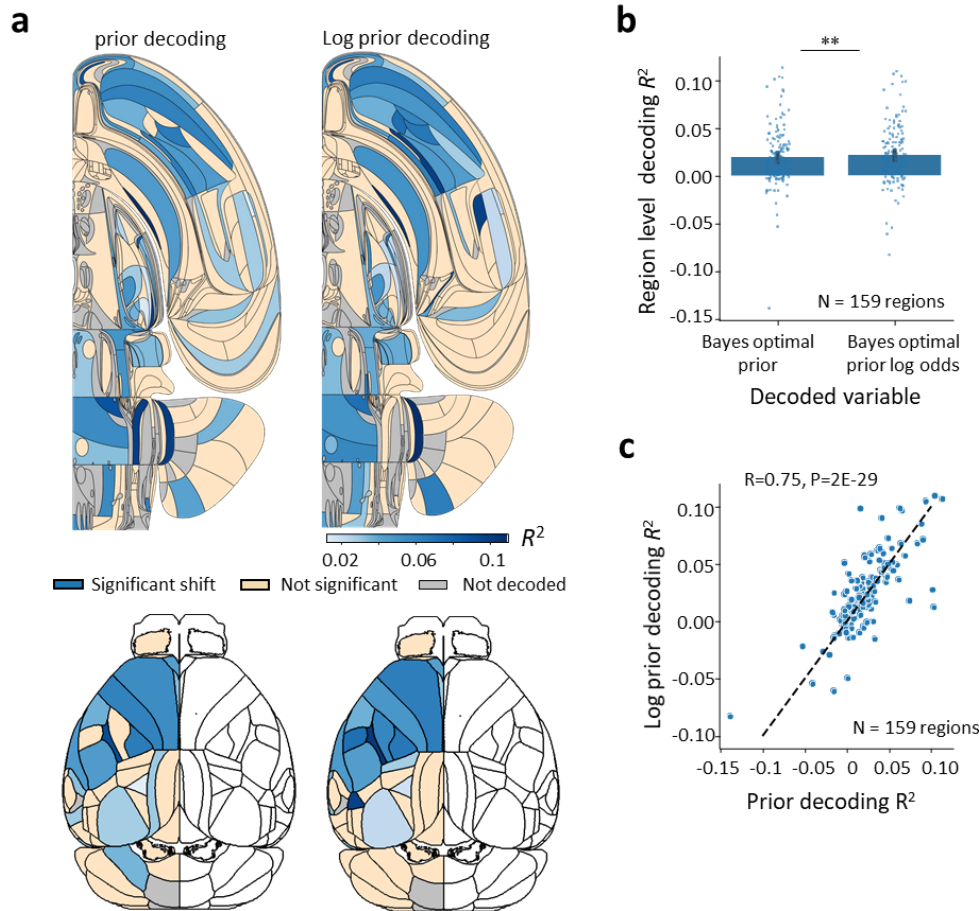


**Figure S2.** Average wheel speed averaged across sessions before and after stimulus onset. The decoded time window used for Ephys data is indicated in light gray. For WFI, the data was decoded on the second-to-last frame relative to the stimulus onset, corresponding to a time window of (-132ms,-66ms)

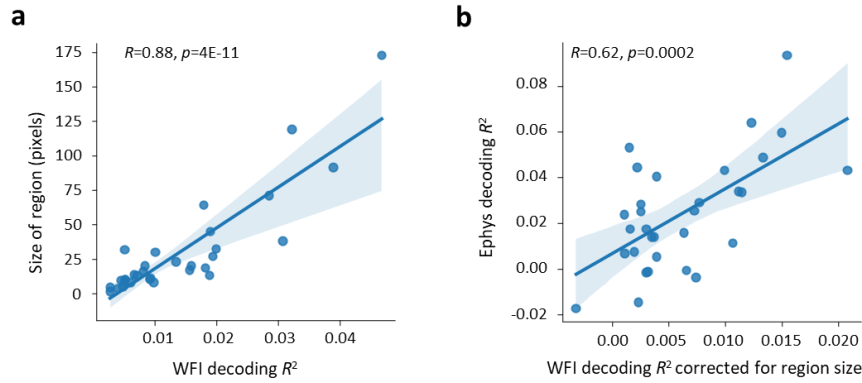


**Figure S3.** Six examples of neurons encoding the Bayes-optimal prior significantly (\*\*\*)  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$ ). Left plots show the spike counts of the neurons (purple line) during the intertrial interval in the [-600, -100] millisecond time window before stimulus onset, along with the Bayes-optimal prior (blue) for a subset of trials within the session. Right subplots correspond to the PSTHs, for all trials within the session, conditioned on the Bayes-optimal prior for the right side being less than 0.3 (orange) vs greater than 0.7 (blue).

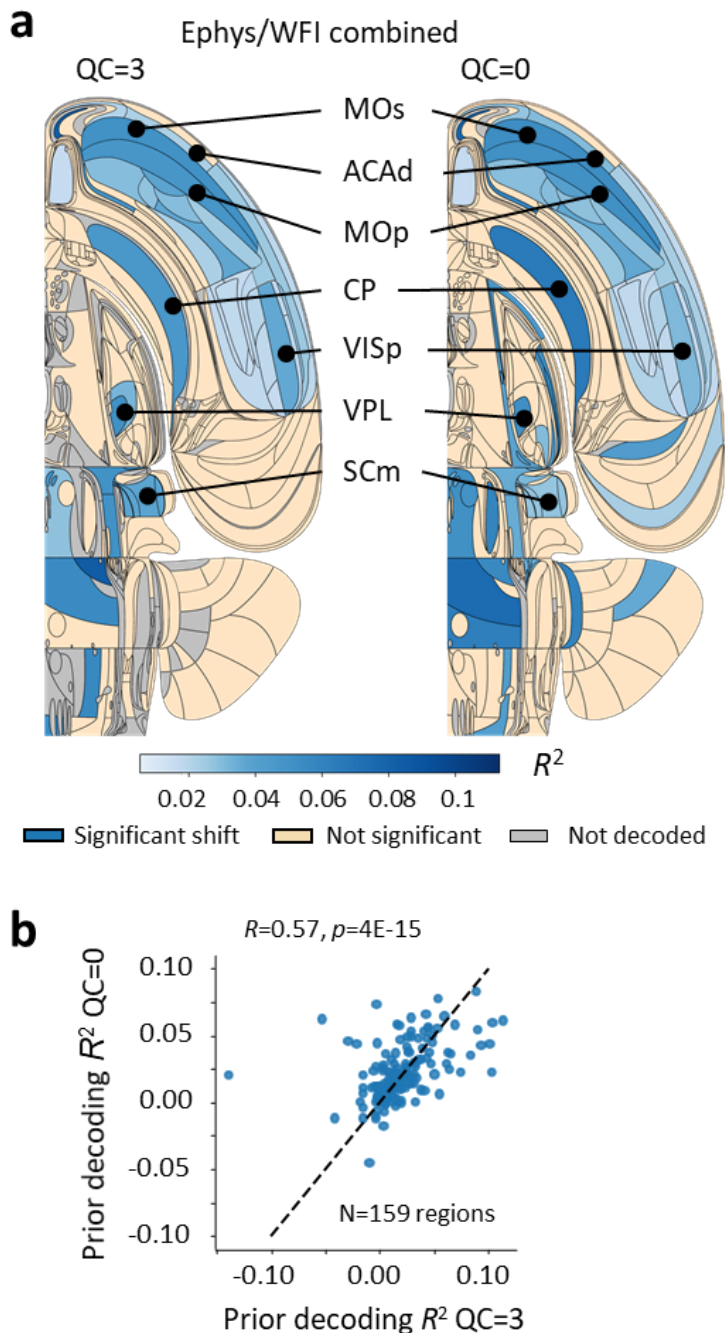




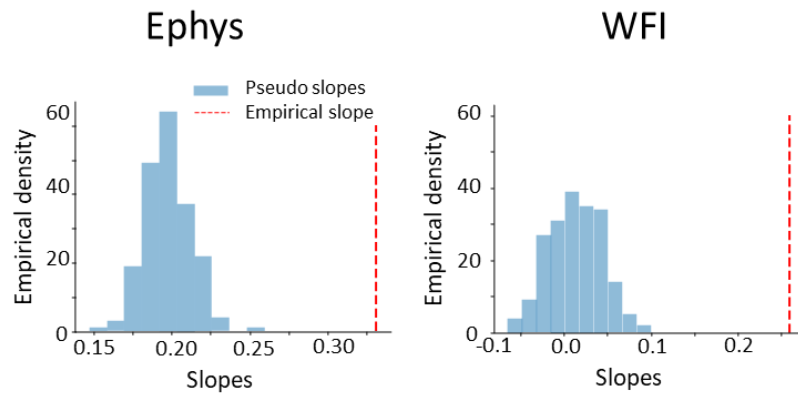
**Figure S4. a.** Swanson maps (top) and dorsal cortex views (bottom) of decoding  $R^2$  for the Bayes-optimal prior (left column) versus the log odds of the Bayes-optimal prior (right column) from Ephys data. These maps are not corrected for multiple comparisons. **b.** The mean decoding  $R^2$  values are significantly higher for the log odds of the Bayes-optimal prior compared to the Bayes-optimal prior (\*,  $p=4.15e-02$ ). **c.** The decoding  $R^2$  of the log odds of the Bayes-optimal prior and the Bayes-optimal prior are highly correlated for Ephys data.



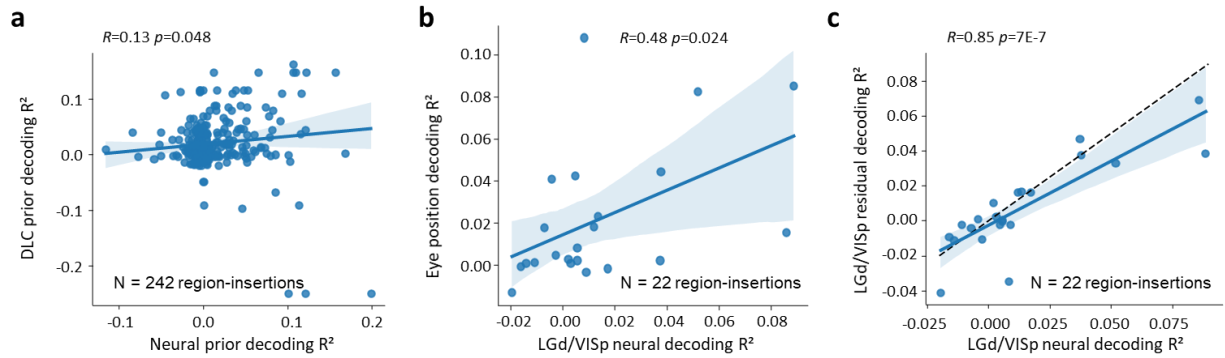
**Figure S5:** **a.** WFI decoding  $R^2$  for the Bayes-optimal prior are significantly correlated with the size of the decoded regions (Spearman correlation  $R=0.88$ ,  $p=4E-11$ ). **b.** The corrected  $R^2$  for Ephys are significantly correlated (Spearman correlation  $R=0.62$ ,  $p=0.0002$ ) even when correcting the WFI  $R^2$  data for region size. Correcting for the region size in WFI was performed by subtracting the size predicted  $R^2$  (from panel a) from the WFI  $R^2$ . Each dot corresponds to one region. All Ephys regions (significant and non-significant) were included in this analysis.



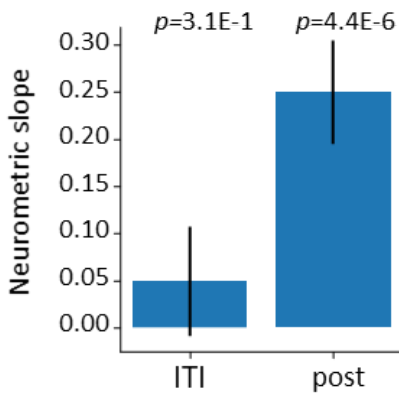
**Figure S6. a.** Ephys/WFI-combined Swanson maps of decoding  $R^2$  for units that have passed all three quality control criteria (QC=3) versus all units returned by the spike sorter (QC=0, same as in Fig. 2b).  $R^2$  are reported for significant regions in blue. Significance is assessed with the Benjamini-Hochberg procedure, correcting for multiple comparisons, with a false discovery rate of 1%. Note that while the number of decoded brain regions is much higher for QC=0 (267 vs 159), the percentage of significant regions (21.3%) remains similar to that obtained with QC=3 (26.8%). **b.** Decoding  $R^2$  are significantly correlated for the 159 regions included in both analyses.



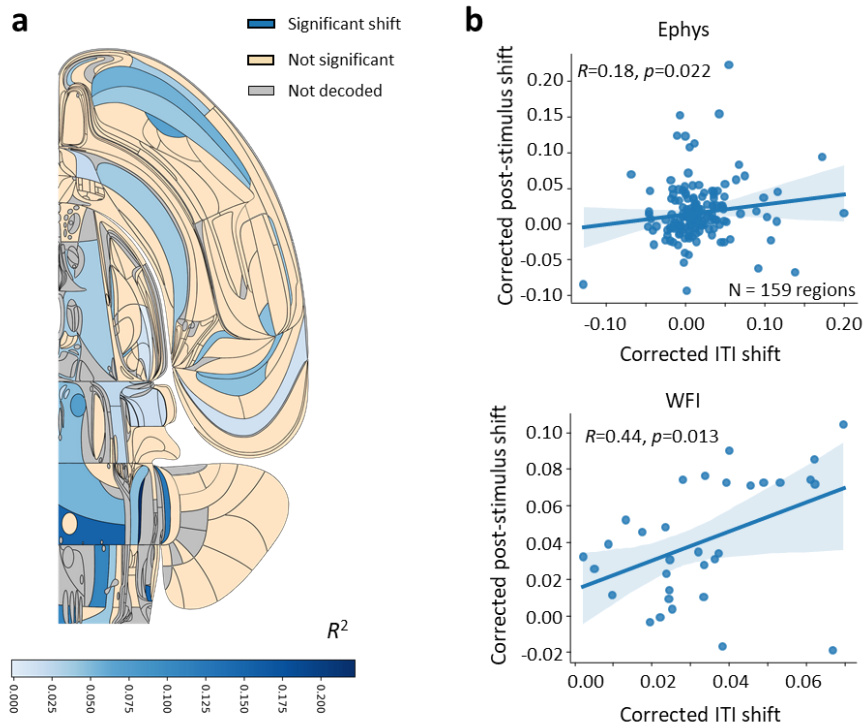
**Figure S7.** Null distribution of the slopes for the proportion of right choice vs decoded prior on zero contrast trials. Slopes were estimated using logistic regression to predict the choice (left or right) as a function of the decoded prior. The null distribution is obtained with  $N=200$  pseudosessions. For each pseudosession, pseudoactions were generated by fitting first each session with an action kernel behavioral model which was then used to generate pseudoactions (see Methods). We can then obtain pseudoslopes by predicting (with logistic regression) the pseudoactions as a function of the decoded prior. The null distribution is obtained by averaging the pseudoslopes across all sessions (we thus obtain  $N=200$  averaged pseudoslopes). The empirical average slope (red dashed lines, values corresponding to Fig. 2e) does not overlap with the null distribution obtained with pseudosessions (blue histogram). Therefore the correlations between the predictions prior and proportion of right choice can not be explained away by spurious temporal correlations or drift in the neural recordings.



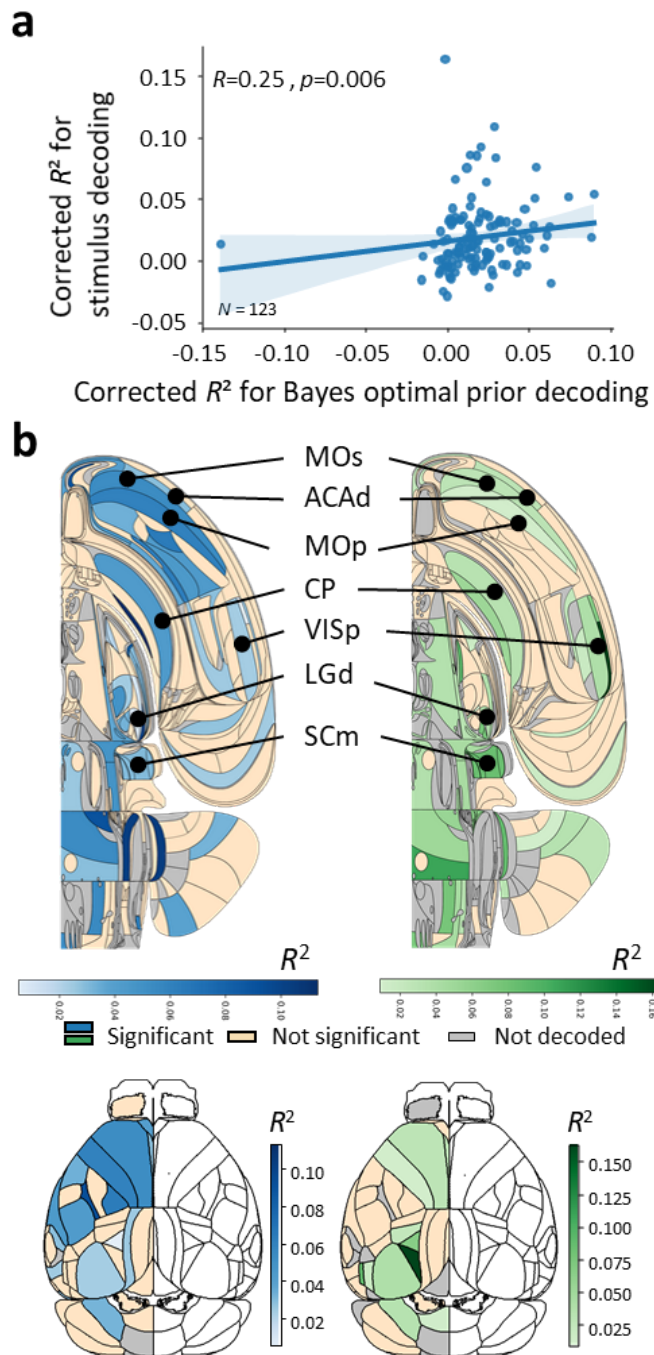
**Figure S8.** **a.** The decoding  $R^2$  for the Bayes-optimal prior from neural activity is significantly correlated with the decoding  $R^2$  for the Bayes-optimal prior from DLC features (Spearman correlation  $R=0.13$ ,  $p=0.048$ ). **b.** Decoding  $R^2$  for the Bayes-optimal prior from neural activity in VISp and LGd against decoding  $R^2$  for the Bayes-optimal prior from eye position. The correlation between these two quantities is significant (Pearson correlation  $R=0.48$ ,  $p=0.024$ ). **c.** Residual decoding  $R^2$  against neural decoding  $R^2$ . The residual decoding  $R^2$  values are obtained by first regressing the Bayes-optimal prior against eye position and then regressing the prior residual (Bayes-optimal prior minus Bayes-optimal prior estimated from eye position) against neural activity in VISp and LGd. The neural decoding  $R^2$  corresponds to the  $R^2$  when decoding the Bayes-optimal prior from neural activity. The two quantities are strongly correlated (Pearson correlation  $R=0.85$ ,  $p=7E-7$ ), suggesting that the prior signals in LGd and VISp are not solely due to the position of the eyes across blocks.



**Figure S9:** The average slope of the neurometric curves is significantly different from 0 during the post stimulus period (2-tailed Wilcoxon paired test,  $t=3689$ ,  $p=4.4E-6$ ,  $N=159$  regions) but not during the ITI ( $t=5768$ ,  $p=0.31$ ,  $N=159$  regions).

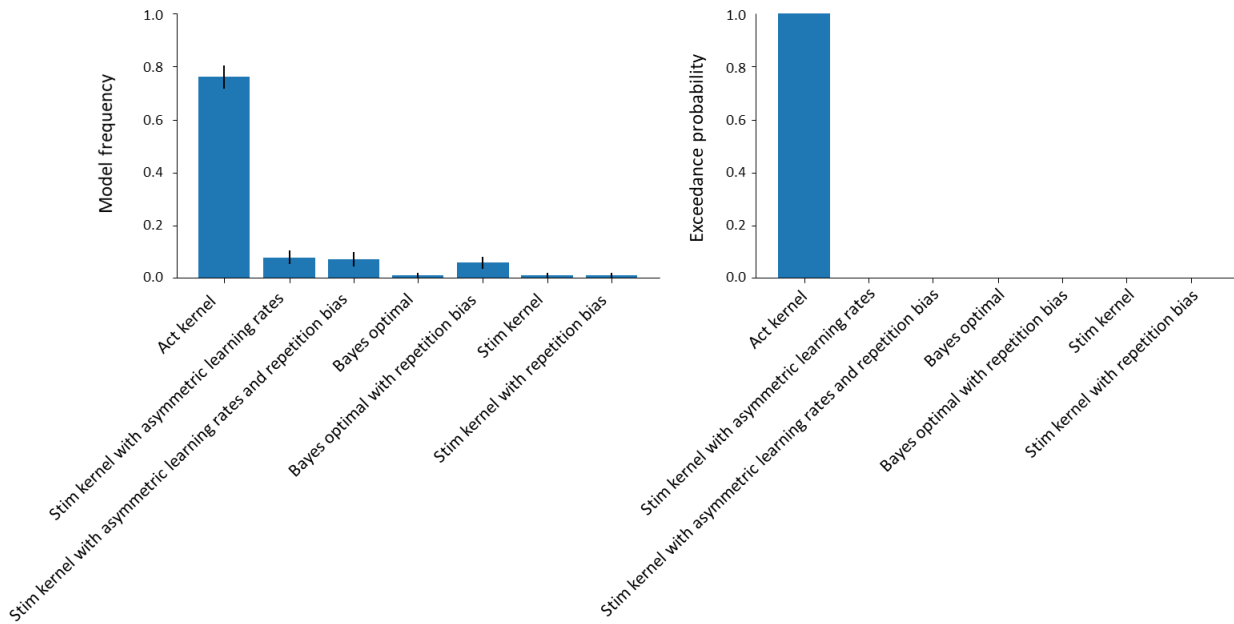


**Figure S10:** a. Swanson map of corrected neurometric posts-stimulus shifts for Ephys data. b. The corrected post-stimulus shifts and corrected ITI shifts are significantly correlated in both Ephys (Spearman correlation  $R=0.18$ ,  $p=0.022$ ,  $N=159$  regions) and WFI (Spearman correlation  $R=0.44$ ,  $p=0.013$ ,  $N=32$  regions).

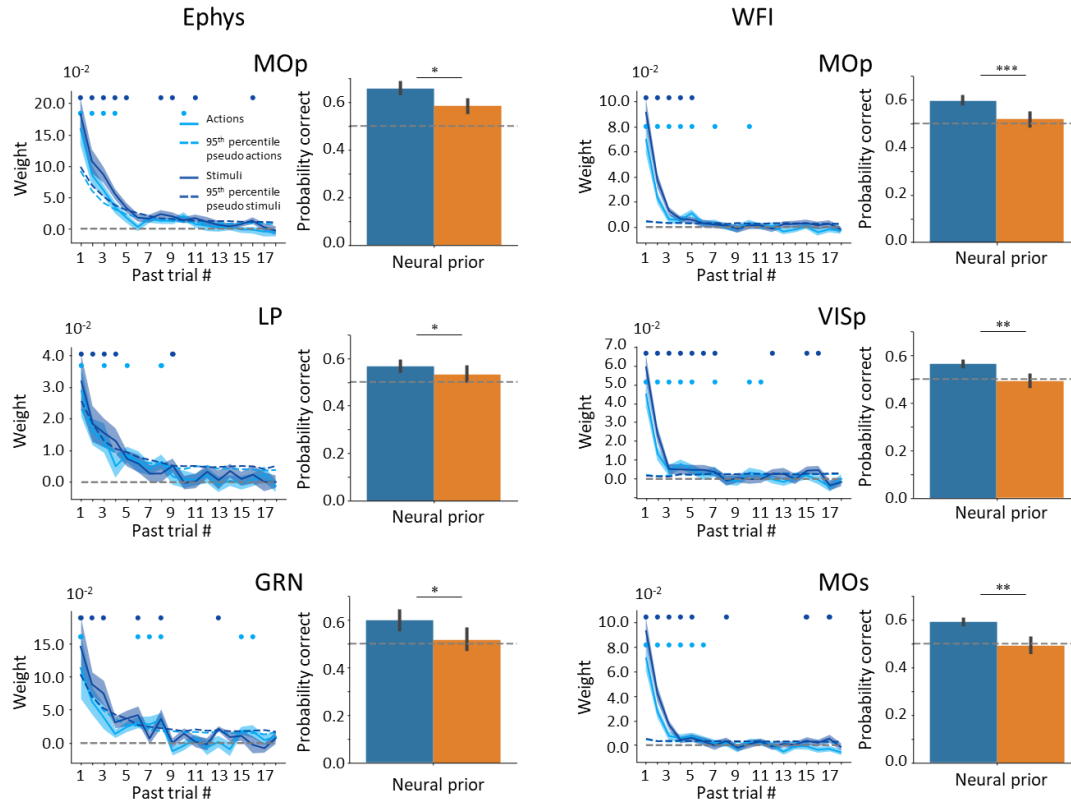


**Figure S11: a.** The neural decoding  $R^2$  for the stimulus and the Bayes-optimal prior are significantly correlated across brain regions (Spearman correlation  $R=0.25, p=0.006$ ). **b.** Swanson maps and dorsal cortical views of brain regions encoding the Bayes-optimal prior (blue, left) and the stimulus (green, right) significantly based on Ephys data.

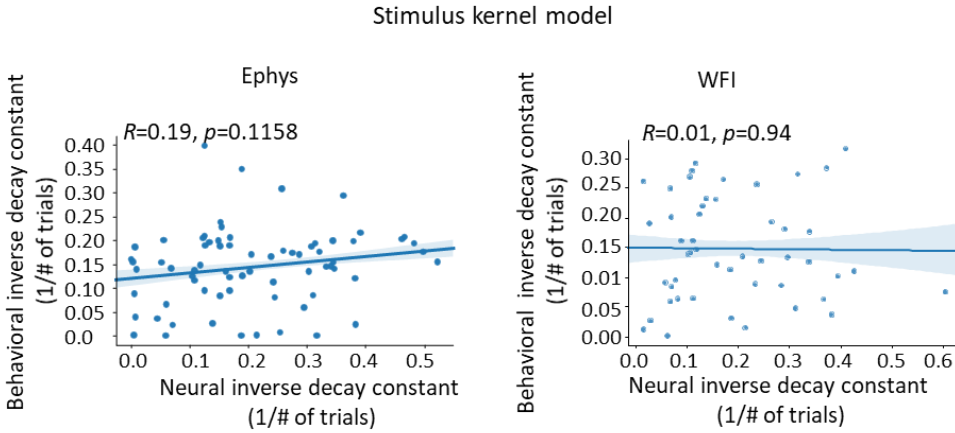




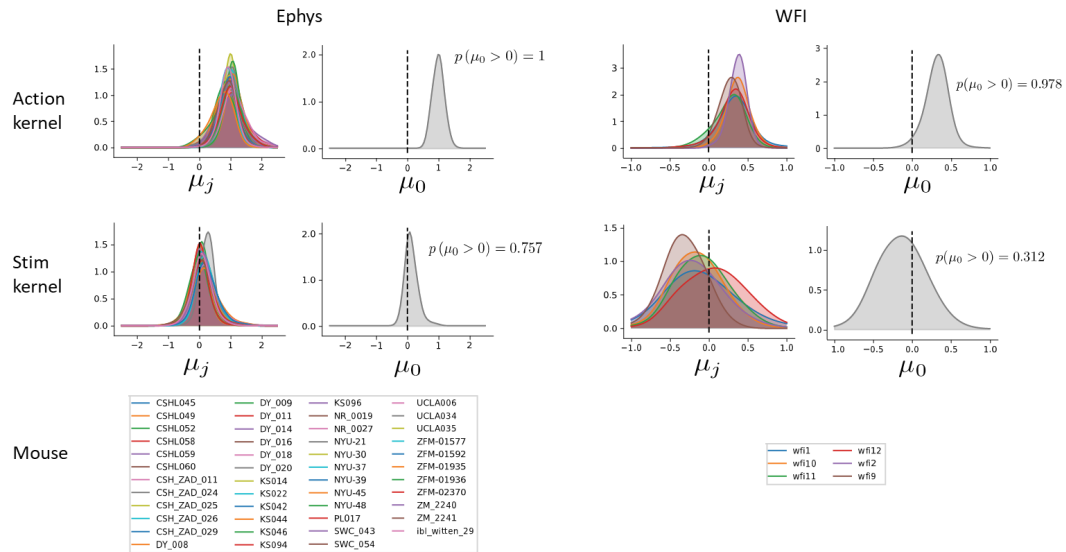
**Figure S12.** Bayesian model comparison for 7 behavioral models, considering the possibility of one step repetition bias (i.e. a tendency to repeat the previous choice), and, for the stimulus kernel model, the presence of positivity and confirmation biases as asymmetric learning rates (Palminteri & Lebreton, 2022). See Methods for more details on the Bayes-optimal, action kernel and stimulus kernel models and Supplementary information for the formal equations of the repetition bias and asymmetrical learning rates. Model frequency (the posterior probability of the model given the subjects' data, left panel) and exceedance probability (the probability that a model is more likely than any other models, right panel) are shown. The action kernel model offered the best account of the data even when including models with repetition, positivity and confirmation biases ( $p_{\text{exceedance}} > 0.999$ ).



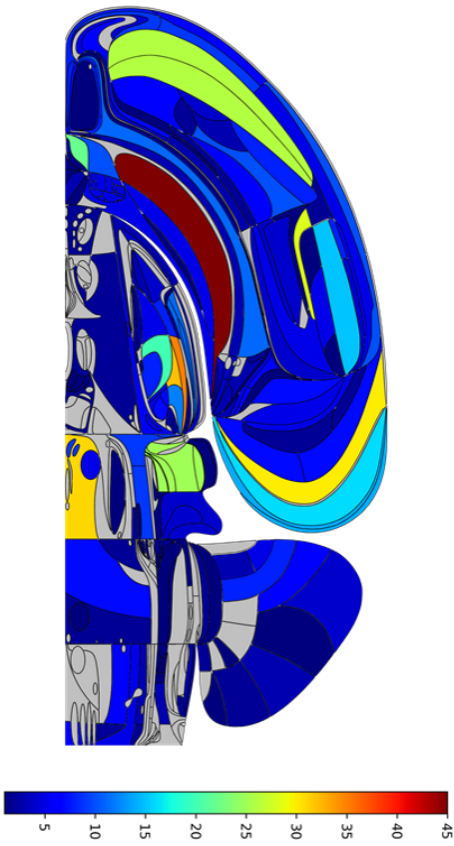
**Figure S13:** Same analysis as in Fig. 4c,e, but for three specific brain regions using Ephys (MOp, LP, GRN) or WFI data (right column, MOp, VISp, MOs) (\*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ ). As the asymmetry effect is observed on a brain-wide level, we performed a 1-tailed signed-rank Wilcoxon paired test for assessing asymmetry significance on a region-level.



**Figure S14.** Behavioral decay constants, obtained by fitting the stimulus kernel model to the behavior, as a function of the neural decay constants, obtained by estimating the temporal dependency of the neural signals with respect to previous stimulus. The neural and behavioral inverse decay constants are not significantly correlated for either Ephys (Pearson correlation  $R=0.19$ ,  $p=0.1158$ ) or WFI (Pearson correlation  $R=0.01$ ,  $p=0.94$ )



**Figure S15.** Hierarchical modeling of the neural and behavioral inverse decay constants (also referred to here as learning rates). The parameter  $\mu_j$ , defined for each mouse  $j$ , is the slope (the multiplicative coefficient) of the linear regression predicting the neural learning rate from  $\mu_0$  the behavioral learning rate (on the sessions of mouse  $j$ ). These parameters  $\mu_j$  are sampled from a common population level prior with mean  $\mu_0$ . The parameter  $\mu_0$ , defined at the population level, characterizes an overall relationship between neural and behavioral learning rates. We found that the relationship between neural and behavioral learning rates is significantly positive for the action kernel model (top row), both in electrophysiology (left column) and in widefield imaging (right column), which is not the case for the stimulus Kernel model (bottom row). Furthermore, when testing the difference in means of the population level parameter  $\mu_0$  between action and stimulus kernels, we found the positive relationship was significantly greater for the action kernel, both in Ephys and in WFI. Significance was assessed by estimating the means of the  $\mu_0$  distributions for the action and stimulus kernels with the BEST Bayesian test (Kruschke, 2013). In both Ephys and WFI, we found that  $p\left(\overline{\mu_0^{actKernel}} > \overline{\mu_0^{stimKernel}}\right) = 1$  with  $\overline{\mu_0^{actKernel}}$  and  $\overline{\mu_0^{stimKernel}}$  the means of the  $\mu_0$  distributions for the action and stimulus kernels, respectively. Regarding the effect sizes, with the same BEST procedure, we find an effect size of 6.27 in Ephys and 2.91 in widefield (effect sizes greater than 1.3 are commonly considered to be very large (Sullivan & Feinn, 2012)). See Supplementary Information for the full specification of the hierarchical generative model.



**Figure S16.** Number of recordings session per brain regions

## References

- Ackley, D. H., Hinton, G. E., & Sejnowski, T. J. (1985). A Learning Algorithm for Boltzmann Machines\*. *Cognitive Science*, 9(1), 147–169.  
[https://doi.org/10.1207/s15516709cog0901\\_7](https://doi.org/10.1207/s15516709cog0901_7)
- Andrieu, C., & Thoms, J. (2008). A tutorial on adaptive MCMC. *Statistics and Computing*, 18(4), 343–373. <https://doi.org/10.1007/s11222-008-9110-y>
- Ashwood, Z. C., Roy, N. A., Stone, I. R., The International Brain Laboratory, Urai, A. E., Churchland, A. K., Pouget, A., & Pillow, J. W. (2022). Mice alternate between discrete strategies during perceptual decision-making. *Nature Neuroscience*, 25(2), 201–212.  
<https://doi.org/10.1038/s41593-021-01007-z>
- Bell, A. H., Summerfield, C., Morin, E. L., Malecek, N. J., & Ungerleider, L. G. (2016). Encoding of Stimulus Probability in Macaque Inferior Temporal Cortex. *Current Biology*, 26(17), 2280–2290. <https://doi.org/10.1016/j.cub.2016.07.007>
- Berkes, P., Orbán, G., Lengyel, M., & Fiser, J. (2011). Spontaneous Cortical Activity Reveals Hallmarks of an Optimal Internal Model of the Environment. *Science*, 331(6013), 83–87.  
<https://doi.org/10.1126/science.1195870>
- Biderman, D., Whiteway, M. R., Hurwitz, C., Greenspan, N., Lee, R. S., Vishnubhotla, A., Warren, R., Pedraja, F., Noone, D., Schartner, M., Huntenburg, J. M., Khanal, A., Meijer, G. T., Noel, J.-P., Pan-Vazquez, A., Socha, K. Z., Urai, A. E., The International Brain Laboratory, Cunningham, J. P., ... Paninski, L. (2023). *Lightning Pose: Improved animal pose estimation via semi-supervised learning, Bayesian ensembling, and cloud-native open-source tools* [Preprint]. Neuroscience. <https://doi.org/10.1101/2023.04.28.538703>
- Bondy, A. G., Haefner, R. M., & Cumming, B. G. (2018). Feedback determines the structure of correlated variability in primary visual cortex. *Nature Neuroscience*, 21(4), 598–606.

<https://doi.org/10.1038/s41593-018-0089-1>

Busse, L., Ayaz, A., Dhruv, N. T., Katzner, S., Saleem, A. B., Schölvinc, M. L., Zaharia, A. D., & Carandini, M. (2011). The Detection of Visual Contrast in the Behaving Mouse. *The Journal of Neuroscience*, 31(31), 11351–11361.

<https://doi.org/10.1523/JNEUROSCI.6689-10.2011>

Dabney, W., Kurth-Nelson, Z., Uchida, N., Starkweather, C. K., Hassabis, D., Munos, R., & Botvinick, M. (2020). A distributional code for value in dopamine-based reinforcement learning. *Nature*, 577(7792), 671–675. <https://doi.org/10.1038/s41586-019-1924-6>

Echeveste, R., Aitchison, L., Hennequin, G., & Lengyel, M. (2020). Cortical-like dynamics in recurrent circuits optimized for sampling-based probabilistic inference. *Nature Neuroscience*, 23(9), 1138–1149. <https://doi.org/10.1038/s41593-020-0671-1>

Elber-Dorozko, L., & Loewenstein, Y. (2018). Striatal action-value neurons reconsidered. *eLife*, 7, e34248. <https://doi.org/10.7554/eLife.34248>

Ernst, M. O., & Banks, M. S. (2002). Humans integrate visual and haptic information in a statistically optimal fashion. *Nature*, 415(6870), 429–433.

<https://doi.org/10.1038/415429a>

Forstmann, B. U. (2010). The neural substrate of prior information in perceptual decision making: A model-based analysis. *Frontiers in Human Neuroscience*, 4.

<https://doi.org/10.3389/fnhum.2010.00040>

Ganguli, D., & Simoncelli, E. P. (2014). Efficient Sensory Encoding and Bayesian Inference with Heterogeneous Neural Populations. *Neural Computation*, 26(10), 2103–2134.

[https://doi.org/10.1162/NECO\\_a\\_00638](https://doi.org/10.1162/NECO_a_00638)

Haefner, R. M., Berkes, P., & Fiser, J. (2016). Perceptual Decision-Making as Probabilistic Inference by Neural Sampling. *Neuron*, 90(3), 649–660.

<https://doi.org/10.1016/j.neuron.2016.03.020>

Han, S., & Helmchen, F. (2023). *Behavior-relevant top-down cross-modal predictions in mouse*

- neocortex* [Preprint]. Neuroscience. <https://doi.org/10.1101/2023.04.03.535389>
- Hanks, T. D., Mazurek, M. E., Kiani, R., Hopp, E., & Shadlen, M. N. (2011). Elapsed Decision Time Affects the Weighting of Prior Probability in a Perceptual Decision Task. *Journal of Neuroscience*, 31(17), 6339–6352. <https://doi.org/10.1523/JNEUROSCI.5613-10.2011>
- Hansen, K. A., Hillenbrand, S. F., & Ungerleider, L. G. (2012). Human Brain Activity Predicts Individual Differences in Prior Knowledge Use during Decisions. *Journal of Cognitive Neuroscience*, 24(6), 1462–1475. [https://doi.org/10.1162/jocn\\_a\\_00224](https://doi.org/10.1162/jocn_a_00224)
- Harris, K. D. (2020). *Nonsense correlations in neuroscience* [Preprint]. Neuroscience. <https://doi.org/10.1101/2020.11.29.402719>
- Hoyer, P. O., & Hyvärinen, P. A. (2003). Interpreting neural response variability as Monte Carlo sampling of the posterior. In *Advances in Neural Information Processing Systems* (Vol. 15, pp. 293–300).
- International Brain Lab, Benson, B., Benson, J., Birman, D., Bonacchi, N., Carandini, M., Catarino, J. A., Chapuis, G. A., Churchland, A. K., Dan, Y., Dayan, P., DeWitt, E. E., Engel, T. A., Fabbri, M., Faulkner, M., Fiete, I. R., Findling, C., Freitas-Silva, L., Gercek, B., ... Witten, I. B. (2023). *A Brain-Wide Map of Neural Activity during Complex Behaviour* [Preprint]. Neuroscience. <https://doi.org/10.1101/2023.07.04.547681>
- Ishizu, K., Nishimoto, S., & Funamizu, A. (2023). *Localized and global computation for integrating prior value and sensory evidence in the mouse cerebral cortex* [Preprint]. Neuroscience. <https://doi.org/10.1101/2023.06.06.543645>
- Jacobs, R. A. (1999). Optimal integration of texture and motion cues to depth. *Vision Research*, 39(21), 3621–3629. [https://doi.org/10.1016/S0042-6989\(99\)00088-7](https://doi.org/10.1016/S0042-6989(99)00088-7)
- Jardri, R., Duverne, S., Litvinova, A. S., & Denève, S. (2017). Experimental evidence for circular inference in schizophrenia. *Nature Communications*, 8(1), 14218. <https://doi.org/10.1038/ncomms14218>
- Knill, D. C., & Pouget, A. (2004). The Bayesian brain: The role of uncertainty in neural coding



and computation. *Trends in Neurosciences*, 27(12), 712–719.

<https://doi.org/10.1016/j.tins.2004.10.007>

Kok, P., Jehee, J. F. M., & de Lange, F. P. (2012). Less Is More: Expectation Sharpens Representations in the Primary Visual Cortex. *Neuron*, 75(2), 265–270.

<https://doi.org/10.1016/j.neuron.2012.04.034>

Krasniak, C. (2022). *Mesoscale imaging, inactivation, and collaboration in a standardized visual decision-making task*. Cold Spring Harbor Laboratory.

<http://repository.cshl.edu/id/eprint/40616/>

Kruschke, J. K. (2013). Bayesian estimation supersedes the t test. *Journal of Experimental Psychology: General*, 142(2), 573–603. <https://doi.org/10.1037/a0029146>

Laboratory, International Brain. (2022). *Video hardware and software for the International Brain Laboratory*. 7030489 Bytes. <https://doi.org/10.6084/M9.FIGSHARE.19694452.V1>

Laboratory, International Brain. (2023). *Data release—Brainwide map—Q4 2022*. 12507703 Bytes. <https://doi.org/10.6084/M9.FIGSHARE.21400815.V6>

Lak, A., Okun, M., Moss, M. M., Gurnani, H., Farrell, K., Wells, M. J., Reddy, C. B., Kepecs, A., Harris, K. D., & Carandini, M. (2020). Dopaminergic and Prefrontal Basis of Learning from Sensory Confidence and Reward Value. *Neuron*, 105(4), 700–711.e6.

<https://doi.org/10.1016/j.neuron.2019.11.018>

Lange, R. D., & Haefner, R. M. (2022). Task-induced neural covariability as a signature of approximate Bayesian learning and inference. *PLOS Computational Biology*, 18(3), e1009557. <https://doi.org/10.1371/journal.pcbi.1009557>

Ma, W. J., Beck, J. M., Latham, P. E., & Pouget, A. (2006). Bayesian inference with probabilistic population codes. *Nature Neuroscience*, 9(11), 1432–1438.

<https://doi.org/10.1038/nn1790>

Mamassian, P., Knill, D. C., & Kersten, D. (1998). The perception of cast shadows. *Trends in Cognitive Sciences*, 2(8), 288–295. [https://doi.org/10.1016/S1364-6613\(98\)01204-2](https://doi.org/10.1016/S1364-6613(98)01204-2)

- Mathis, A., Mamidanna, P., Cury, K. M., Abe, T., Murthy, V. N., Mathis, M. W., & Bethge, M. (2018). DeepLabCut: Markerless pose estimation of user-defined body parts with deep learning. *Nature Neuroscience*, *21*(9), 1281–1289.  
<https://doi.org/10.1038/s41593-018-0209-y>
- Mayrhofer, J. M., El-Boustani, S., Foustoukos, G., Auffret, M., Tamura, K., & Petersen, C. C. H. (2019). Distinct Contributions of Whisker Sensory Cortex and Tongue-Jaw Motor Cortex in a Goal-Directed Sensorimotor Transformation. *Neuron*, *103*(6), 1034-1043.e5.  
<https://doi.org/10.1016/j.neuron.2019.07.008>
- Mendonça, A. G., Drugowitsch, J., Vicente, M. I., DeWitt, E. E. J., Pouget, A., & Mainen, Z. F. (2020). The impact of learning on perceptual decisions and its implication for speed-accuracy tradeoffs. *Nature Communications*, *11*(1), 2757.  
<https://doi.org/10.1038/s41467-020-16196-7>
- Mulder, M. J., Wagenmakers, E.-J., Ratcliff, R., Boekel, W., & Forstmann, B. U. (2012). Bias in the Brain: A Diffusion Model Analysis of Prior Probability and Potential Payoff. *Journal of Neuroscience*, *32*(7), 2335–2343. <https://doi.org/10.1523/JNEUROSCI.4156-11.2012>
- Niv, Y. (2019). Learning task-state representations. *Nature Neuroscience*, *22*(10), 1544–1553.  
<https://doi.org/10.1038/s41593-019-0470-8>
- Nogueira, R., Abolafia, J. M., Drugowitsch, J., Balaguer-Ballester, E., Sanchez-Vives, M. V., & Moreno-Bote, R. (2017). Lateral orbitofrontal cortex anticipates choices and integrates prior with current information. *Nature Communications*, *8*(1), 14823.  
<https://doi.org/10.1038/ncomms14823>
- Norton, E. H., Acerbi, L., Ma, W. J., & Landy, M. S. (2019). Human online adaptation to changes in prior probability. *PLOS Computational Biology*, *15*(7), e1006681.  
<https://doi.org/10.1371/journal.pcbi.1006681>
- Palminteri, S., & Lebreton, M. (2022). The computational roots of positivity and confirmation biases in reinforcement learning. *Trends in Cognitive Sciences*, *26*(7), 607–621.

<https://doi.org/10.1016/j.tics.2022.04.005>

- Park, J., Kim, S., Kim, H. R., & Lee, J. (2022). *Prior expectation enhances sensorimotor behavior by modulating population tuning and subspace activity in the sensory cortex* [Preprint]. *Neuroscience*. <https://doi.org/10.1101/2022.12.04.516847>
- Pedregosa, F., Varoquaux, G., Gramfort, A., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Coupaneau, D., Brucher, M., Perrot, M., & Duchesnay, E. (2011). *Scikit-learn: Machine Learning in Python*. *12*, 2825–2830.
- Platt, M. L., & Glimcher, P. W. (1999). Neural correlates of decision variables in parietal cortex. *Nature*, *400*(6741), 233–238. <https://doi.org/10.1038/22268>
- Rao, V., DeAngelis, G. C., & Snyder, L. H. (2012). Neural Correlates of Prior Expectations of Motion in the Lateral Intraparietal and Middle Temporal Areas. *Journal of Neuroscience*, *32*(29), 10063–10074. <https://doi.org/10.1523/JNEUROSCI.5948-11.2012>
- Sahani, M., & Dayan, P. (2003). Doubly Distributional Population Codes: Simultaneous Representation of Uncertainty and Multiplicity. *Neural Computation*, *15*(10), 2255–2279. <https://doi.org/10.1162/089976603322362356>
- Schaeffer, R., Khona, M., Meshulam, L., International Brain Laboratory, & Fiete, I. R. (2020). *Reverse-engineering Recurrent Neural Network solutions to a hierarchical inference task for mice* [Preprint]. *Neuroscience*. <https://doi.org/10.1101/2020.06.09.142745>
- Scott, S. L. (2002). Bayesian Methods for Hidden Markov Models: Recursive Computing in the 21st Century. *Journal of the American Statistical Association*, *97*(457), 337–351. <https://doi.org/10.1198/016214502753479464>
- Soltani, A., & Wang, X.-J. (2010). Synaptic computation underlying probabilistic inference. *Nature Neuroscience*, *13*(1), 112–119. <https://doi.org/10.1038/nn.2450>
- Son, S., Moon, J., Kim, Y.-J., Kang, M.-S., & Lee, J. (2023). Frontal-to-visual information flow explains predictive motion tracking. *NeuroImage*, *269*, 119914. <https://doi.org/10.1016/j.neuroimage.2023.119914>

- Stephan, K. E., Penny, W. D., Daunizeau, J., Moran, R. J., & Friston, K. J. (2009). Bayesian model selection for group studies. *NeuroImage*, *46*(4), 1004–1017.  
<https://doi.org/10.1016/j.neuroimage.2009.03.025>
- Sugawara, M., & Katahira, K. (2021). Dissociation between asymmetric value updating and perseverance in human reinforcement learning. *Scientific Reports*, *11*(1), 3574.  
<https://doi.org/10.1038/s41598-020-80593-7>
- Sullivan, G. M., & Feinn, R. (2012). Using Effect Size—Or Why the *P* Value Is Not Enough. *Journal of Graduate Medical Education*, *4*(3), 279–282.  
<https://doi.org/10.4300/JGME-D-12-00156.1>
- The International Brain Laboratory, Aguillon-Rodriguez, V., Angelaki, D., Bayer, H., Bonacchi, N., Carandini, M., Cazettes, F., Chapuis, G., Churchland, A. K., Dan, Y., Dewitt, E., Faulkner, M., Forrest, H., Haetzel, L., Häusser, M., Hofer, S. B., Hu, F., Khanal, A., Krasniak, C., ... Zador, A. M. (2021). Standardized and reproducible measurement of decision-making in mice. *ELife*, *10*, e63711. <https://doi.org/10.7554/eLife.63711>
- Walker, E. Y., Cotton, R. J., Ma, W. J., & Tolias, A. S. (2020). A neural basis of probabilistic computation in visual cortex. *Nature Neuroscience*, *23*(1), 122–129.  
<https://doi.org/10.1038/s41593-019-0554-5>
- Wang, Q., Ding, S.-L., Li, Y., Royall, J., Feng, D., Lesnar, P., Graddis, N., Naeemi, M., Facer, B., Ho, A., Dolbeare, T., Blanchard, B., Dee, N., Wakeman, W., Hirokawa, K. E., Szafer, A., Sunkin, S. M., Oh, S. W., Bernard, A., ... Ng, L. (2020). The Allen Mouse Brain Common Coordinate Framework: A 3D Reference Atlas. *Cell*, *181*(4), 936–953.e20.  
<https://doi.org/10.1016/j.cell.2020.04.007>
- Weiss, Y., Simoncelli, E. P., & Adelson, E. H. (2002). Motion illusions as optimal percepts. *Nature Neuroscience*, *5*(6), 598–604. <https://doi.org/10.1038/nn0602-858>
- Zemel, R. S., Dayan, P., & Pouget, A. (1998). Probabilistic Interpretation of Population Codes. *Neural Computation*, *10*(2), 403–430. <https://doi.org/10.1162/089976698300017818>