

# Task learning reveals signatures of sample-based internal models in rodent prefrontal cortex

Abhinav Singh<sup>a</sup>, Adrien Peyrache<sup>b</sup>, and Mark D Humphries<sup>a,1</sup>

<sup>a</sup>Faculty of Biology, Medicine and Health, University of Manchester, Manchester, M13 9PT, UK; <sup>b</sup>McGill University, Montreal, Canada

This manuscript was compiled on July 14, 2016

**The inherent uncertainty of the world suggests that brains use probabilistic internal models. We sought to test whether or how neurons represent such models in higher cortical regions, learn them, and use them in behaviour. Using a sampling framework, we predicted that trial-evoked and sleeping population activity represent the inferred and expected probabilities generated from an internal model of a behavioural task, and would become more similar as the task was learnt. To test these predictions, we analysed population activity from rodent prefrontal cortex before, during, and after sessions of learning rules on a maze. Distributions of activity patterns converged between trials and post-learning sleep during successful rule learning. Learning induced changes were greatest for patterns predicting correct choice and expressed at the maze's choice point, consistent with an updated internal model of the task. Our results suggest sample-based internal models are a general computational principle of cortex.**

neural ensembles | statistical inference | Bayes theorem

How do we know what state the world is in? Behavioural evidence suggests brains solve this problem using probabilistic reasoning [1, 2]. Such reasoning implies that brains represent and learn internal models for the statistical structure of the external world [1, 3, 4]. With these models, neurons could represent uncertainty about the world with probability distributions, and update those distributions with new knowledge using the rules of probabilistic inference. Theoretical work has elucidated potential mechanisms for how cortical populations represent and compute with probabilities [2, 5–8], and shown how computational models of inference predict aspects of cortical activity in sensory and decision-making tasks [e.g. 2, 9]. But experimental evidence for the neural representation of probabilistic internal models is lacking.

An experimentally-accessible proposal is the recent inference-by-sampling hypothesis [7, 10–12]. This proposes that cortical population activity at some time  $t$  is a sample from an underlying probability distribution. Cortical activity evoked by external stimuli represents sampling from the model-generated “posterior” distribution that the world is in a particular state. Spontaneous cortical activity represents sampling of the model in the absence of external stimuli, forming a model-generated “prior” for the expected properties of the world. A key prediction is that the evoked and spontaneous population activity should converge over repeated experience, as the internal model adapts to match the relevant statistics of the external world. Just such a convergence has been observed in small populations from ferret V1 over development [11]. Unknown is whether probabilistic internal models are a general computational principle for cortex: whether they can be observed during learning, or in higher-order cortices, or during ongoing behaviour.

A natural candidate to address these issues is the medial

prefrontal cortex (mPFC). Medial PFC is necessary for learning new rules or strategies [13, 14], and changes in mPFC neuron firing times correlates with successful rule learning [15], suggesting that mPFC coding of task-related variables changes over learning. We thus hypothesised that mPFC encodes an internal model of a task, which is updated by task performance, and from which the statistical distributions of population activity are generated. To test these hypotheses, we analysed previously-recorded population activity from the medial prefrontal cortex of rats learning rules in a Y-maze [16].

## Results

Rats with implanted tetrodes learnt one of three rules on a Y-maze: go left, go right, or go to the randomly-lit arm (Fig. 1A). Each recording session was a single day containing 3 epochs totalling typically 1.5 hours: pre-task sleep/rest, behavioural testing on the task, and post-task sleep/rest. We focussed on ten sessions where the animal reached the learning criteria for a rule mid-session (Materials and Methods; 15–55 neurons per session). In this way, we sought to isolate changes in population activity solely due to rule-learning.

**Theory sketch.** Here we outline our theoretical predictions for changes in population activity, derived from the inference-by-sampling hypothesis; a full account is given in the SI Text. We sought to test the idea that the mPFC contains at least one internal model related to task performance, such as representing the relevant decision-variable (here, left or right) or the rule-dependent outcomes. Learning of the task should therefore

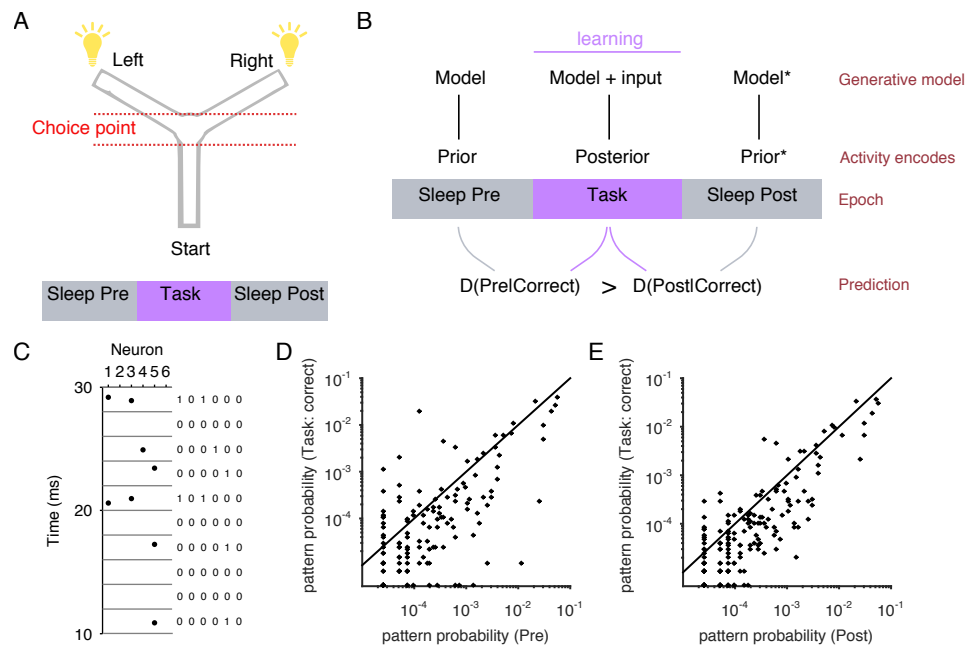
## Significance Statement

The cerebral cortex contains billions of neurons. The activity of one neuron is lost in this morass, so it is thought that the co-ordinated activity of groups of neurons – “neural ensembles” – are the basic element of cortical computation, underpinning sensation, cognition, and action. But what do these ensembles represent? Here we show that ensemble activity in rodent prefrontal cortex represents samples from an internal model of the world - a probability distribution that the world is in a specific state. We find that this internal model is updated during learning about changes to the world, and is sampled during sleep. These results suggest that probability-based computation is a generic principle of cortex.

A.S: designed the analyses; analysed the data; discussed the results; commented on drafts. A.P: discussed the results; commented on drafts. M.D.H: designed the analyses; analysed the data; discussed the results; wrote the paper.

The authors declare no conflict of interest

<sup>1</sup>To whom correspondence should be addressed. E-mail: mark.humphriesmanchester.ac.uk



**Fig. 1.** Activity pattern distributions during rule-learning. (A) Y-maze task set-up (top); each session included the epochs of pre-task sleep/rest, task trials, and post-task sleep/rest (bottom) - Fig. 4A gives a breakdown per session. One of three target rules for obtaining reward was enforced throughout a session: go right; go left; go to the randomly-lit arm. (B) Schematic of theory. If prefrontal cortex encodes an internal model of the task, then activity during the task is derived from the internal model plus the relevant external inputs: the distribution of activity is thus the posterior distribution over the encoded task variables. During sleep, the distribution of activity is derived entirely from the internal model, and thus is the prior distribution over the encoded task variables. Updates to the internal model by task learning (creating Model\*) will then change the prior distribution encoded during sleep (to Prior\*). The theoretical prediction is then that the activity distribution in post-task sleep, derived from the model of the correct rule, will be closer to the distribution on the correct trials, compared to the pre-task sleep. (C) The population activity of simultaneously recorded spike trains was represented as a binary activity pattern in some small time-bin (here 2 ms). (D) Scatter plot of the joint frequency of every occurring pattern in pre-task SWS (distribution  $P(Pre)$ ) and task (distribution  $P(R)$ ) epochs for one session. (E) For the same session as D, scatter plot of the joint frequency of every occurring pattern in post-task SWS [ $P(Post)$ ] and task [ $P(R)$ ] epochs.

update the internal model based on feedback from each trial's outcome. We theorised that mPFC population activity on each trial was sampling from the posterior distribution generated from this model; and that "spontaneous" activity in slow-wave sleep (SWS), occurring in the absence of task-related stimuli and behaviour, samples the corresponding prior distribution (Fig. 1B). Consequently, updating the internal model from task feedback should be reflected in changes to the posterior and prior distributions generated from that model.

By restricting our analyses to sessions with successful learning, we expected the post-task SWS activity to be sampling from an internal model that has learnt the correct rule. To compare posterior distribution samples from the same internal model, we considered population activity during correct trials after the learning criteria were met – we call this distribution  $P(R)$ . Our main prediction was thus that the distribution  $P(R)$  of activity during the correct trials would be more similar to the distribution in post-task SWS [ $P(Post)$ ] than in pre-task SWS [ $P(Pre)$ ]. Such a convergence of distributions would be evidence that a task-related internal model in mPFC was updated by feedback.

**Activity distributions converge between task and post-task sleep.** To test these hypotheses, we compared the statistical distributions of activity patterns between task and sleep epochs. Activity patterns, the putative samples from the underlying probability distribution, were characterised as a binary vector (or "word") of active and inactive neurons with a binsize of 2 ms (Fig. 1C). Each recorded population of  $N$  neurons had the same sub-set of all  $2^N$  possible activity patterns in all epochs (Fig. S1) [17, 18]. Such a common set of patterns is consistent with their being samples generated from the same form of internal model across both behaviour and sleep.

If activity patterns are samples from a probability distribution, then two similar probability distributions will be revealed by the similar frequencies of sampling each pattern [11]. For each pair of epochs, we thus computed the distances between the two corresponding distributions of activity patterns (Fig. 1D,E). We first used the information-theory based Kullback-Liebler divergence to measure the distance  $D(P|Q)$  between distributions  $P$  and  $Q$  in bits [11]. We found that in 9 of the 10 sessions the distribution  $P(R)$  of activity during the trials was closer to the distribution in post-task SWS [ $P(Post)$ ] than in pre-task SWS [ $P(Pre)$ ] (Fig. 2A).

On average the task-evoked distribution of patterns was  $18.7 \pm 6.2\%$  closer to the post-task SWS distribution than the pre-task SWS distribution (Fig. 2B), showing a convergence between task-evoked and post-task SWS distributions. Further, we found a robust convergence even at the level of individual sessions (Fig. 2C).

While the Kullback-Liebler divergence provides the most complete characterisation of the distance between two probability distributions, estimating it accurately from limited sample data has known issues [19]. To check our results were robust, we re-computed all distances using the Hellinger distance, a non-parametric measure that provides a lower bound for the Kullback-Liebler divergence. Reassuringly, we found the same results: the distribution  $P(R)$  of activity during the trials was consistently closer to the distribution in post-task SWS [ $P(Post)$ ] than in pre-task SWS [ $P(Pre)$ ] (Fig. 2F-H; the mean convergence between task-evoked and post-task SWS distributions was  $21 \pm 2.8\%$ ).

The convergence between the task  $P(R)$  and post-task SWS  $P(Post)$  distributions was also robust to both the choice of activity pattern binsize across an order of magnitude from 2 to 20 ms (Fig. S2) and the choice of correct trials in the task distribution  $P(R)$  (Fig. S3).

Together, these results are consistent with the convergence over learning of the posterior and prior distributions represented by mPFC population activity. They imply that mPFC encodes a task-related internal model that is updated by task feedback.

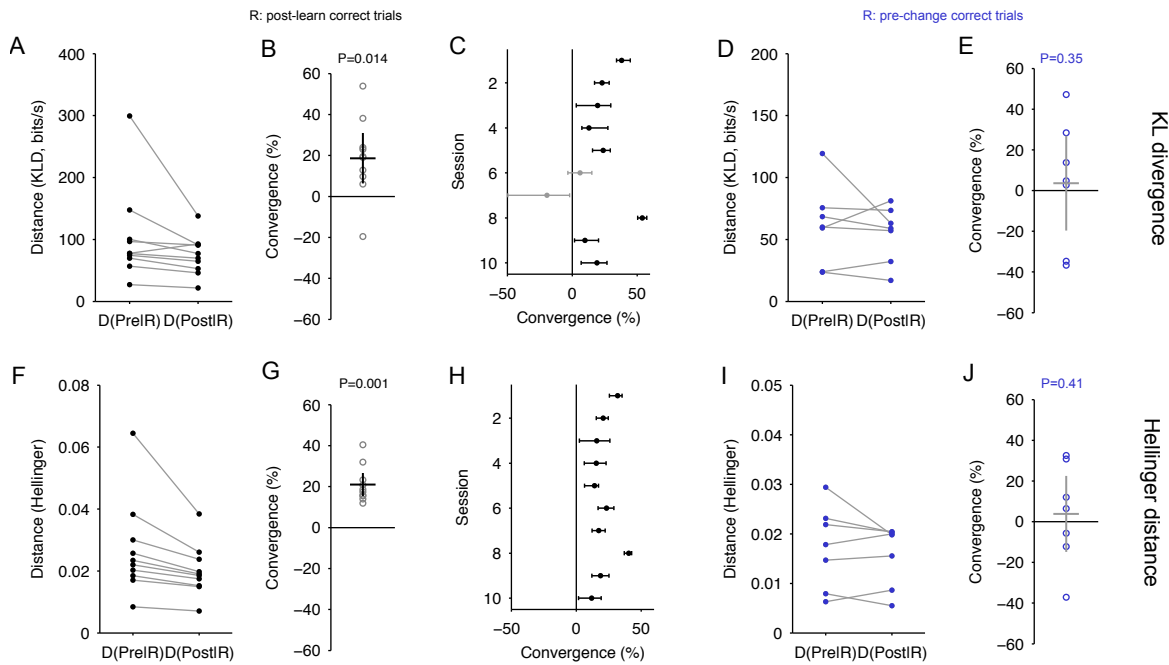
**Wrong models lead to no convergence.** Is this convergence of trial-evoked and post-task SWS distributions inevitable? To answer this, we made use of the 7 sessions in which the rats experienced a rule change. As rule changes occurred only after 10 consecutive correct trials [16], these sessions are uniquely divided into periods when the internal model of the task was right and when it was wrong. Once wrong, the rat needed to find the correct new model. Consequently, our theory predicts that the prior distribution in post-task SWS sleep is not derived from the same internal model as that used before the rule change. In other words, the posterior from the pre-change trials and the prior from the spontaneous activity of post-task sleep are derived from different models, and should not converge.

We tested this prediction by comparing the activity pattern distributions in pre-change correct trials [ $P(R)$ ] and in post-task SWS [ $P(Post)$ ]. There was no convergence between the two distributions, when measured using either Kullback-Liebler divergence ( $3.6 \pm 11.7\%$ ; Fig. 2D-E) or the Hellinger distance ( $3.8 \pm 9.3\%$ ; Fig. 2I-J). For the effect sizes observed for the learning sessions, there was sufficient power to recover the same effect size at  $\alpha = 0.05$  with  $N = 7$  sessions (KLD: learning session effect size  $d = 0.96$ , rule-change session power = 0.7; Hellinger:  $d = 2.36$ , power  $\approx 1$ ), which argues against low power causing the lack of convergence.

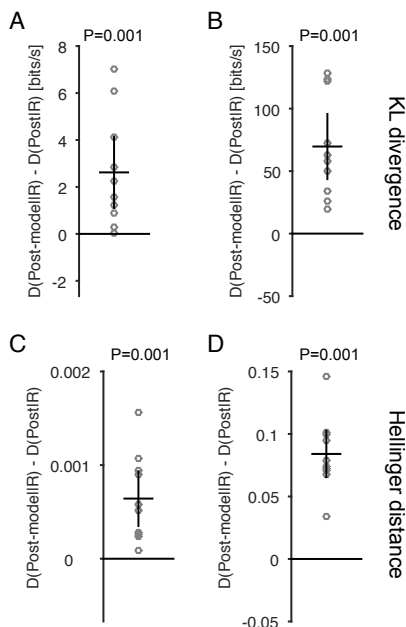
**Convergence is a consequence of changes to correlations, not just firing rates.** Population firing rate differences between waking and sleep states [20, 21], and increases in SWS firing after task-learning, could potentially account for the convergence of distributions during learning. To control for this, we used the "raster" model [20] to generate surrogate sets of spike-trains that matched both the mean firing rates of each neuron, and the distribution of total population activity in each time-bin ( $K = 0, 1, \dots, N$  spikes per bin). Consequently, the occurrence rates of particular activity patterns in the raster model are those predicted to arise from neuron and population firing rates alone.

We found that firing rates could not account for the convergence between task and post-task SWS distributions. The data-derived distance  $D(Post|R)$  was always smaller than the distance  $D(Post - model|R)$  predicted by the raster model (Fig. 3A). This was true whether we used Kullback-Liebler divergence or the Hellinger distance (Fig. 3C) to measure distances between distributions.

Our activity patterns are built from single units, unlike previous work using multi-unit activity [11, 20, 22–24], so we expect our patterns to be sparse with rare synchronous activity. Indeed our data are dominated by activity patterns with  $K = 0$  and  $K = 1$  spikes (Fig. S4). If all patterns were  $K = 0$  or  $K = 1$ , the raster model spike trains would be



**Fig. 2.** Convergence of activity pattern distributions between the task and post-task sleep. (A) Distances between the distributions of pattern frequencies in sleep and task epochs; one dot per session.  $D(X|Y)$ : distance between pattern distributions in epochs  $X$  and  $Y$ : Pre: pre-task SWS; Post: post-task SWS; R: correct task trials. (B) Scatter plot of convergence across all sessions (circles). Convergence is  $D(Pre|R) - D(Post|R)/D(Pre|R)$ . A value greater than zero means that the activity pattern distribution in the task is closer to the distribution in post-task SWS than the distribution in pre-task SWS. Black lines give mean  $\pm$  2 s.e.m. (C) Data (dot) and 95% bootstrapped confidence interval (line) for the convergence of task and post-task SWS activity pattern distributions for each session. Black: sessions with CIs above 0. (D) Distances between the distributions of pattern frequencies in sleep and task epochs during rule-change sessions; one dot per session. Here the distribution  $P(R)$  is constructed from correct task trials before the rule change. (E) Result from panel D expressed as convergence. (F)-(J) As A-E, using Hellinger distance. All  $P$ -values from 1-tailed Wilcoxon signrank test, with  $N=10$  learning sessions (panels A-C and F-H) or  $N = 7$  rule-change sessions (panels D-E and I-J).



**Fig. 3.** Convergence is caused by changes in correlation, not firing rate. (A) The distance between the task and post-task sleep distributions  $D(Post|R)$  is always smaller than predicted by firing rate changes during sleep alone  $D(Post - model|R)$ , as given by the raster model. Black lines give mean  $\pm$  2 s.e.m in all panels. (B) As in A, using only activity patterns with  $K \geq 2$  spikes from data and model. (C)-(D) As A-B, using Hellinger distance. All  $P$ -values from 1-tailed Wilcoxon signrank test, with  $N=10$  sessions.

exactly equivalent to the data. Given the relative sparsity of  $K \geq 2$  patterns in our data, it is all the more surprising then that we found such a consistent lower distance for our data-derived distributions.

It follows that the true difference between data and model is in the relative occurrence of co-activation patterns with  $K \geq 2$  spikes. To check this, we applied the same analysis to distributions built only from these co-activation patterns, drawn from data and from the raster model fitted to the complete data. We found that the data-derived distance  $D(Post|R)$  was always smaller than the distance  $D(Post - model|R)$  predicted by the raster model (Fig. 3B-D). Across all sessions, the model-predicted distance  $D(Post - model|R)$  was between 3% and 46% greater than the data-derived distance  $D(Post|R)$  using Kullback-Liebler divergence, indicating that much of the convergence between task and SWS distributions could not be accounted for by firing rates alone. Consequently, the convergence of task and post-task SWS distributions is not due to population firing rates, but to similarly precise correlations of neuron firing.

Reassuringly, for these  $K \geq 2$  activity pattern distributions, all convergence results held (Fig. S5) despite the order-of-magnitude fewer sampled patterns.

**Convergence is not a recency effect.** We examined periods of SWS in order to most likely observe the sampling of a putative internal model in a static condition, with no external inputs and minimal learning. But as correct task trials more likely occur towards the end of a learning session, this raises the possibility that the closer match between task and post-task

SWS distributions are a recency effect, due to some trace or reverberation in sleep of the most recent task activity.

The time-scales involved make this unlikely. Bouts of SWS did not start until typically 8 minutes after the end of the task (mean 397s, S.D. 188 s; Fig. 4A). Any reverberation would thus have to last at least that long to appear in the majority of post-task SWS distributions.

The intervening period before the first bout of SWS contains quiet wakefulness and early sleep stages. If convergence was a recency effect, then we would expect that distributions [ $P(Rest)$ ] of activity patterns in this more-immediate “rest” epoch would also converge with the task distributions. We did not find this: across sessions, there was no evidence that the distribution in post-task rest [ $P(Rest)$ ] consistently converged on the distribution during task trials [ $P(R)$ ] (Fig. 4B-C; mean convergence was  $-8.7 \pm 18.7\%$ ). Not only is the observed convergence inconsistent with a recency effect, it seems also selective for activity in SWS.

### Distributions are updated by task-relevant activity patterns.

The above analysis rests on the idea that the distributions of activity patterns are derived from an internal model of the task. This predicts that individual patterns should correlate with some aspect of the task. We sought an unbiased way of testing this prediction, so considered the following. In our theory, the changes to the internal model over learning should be directly reflected in the differences between the prior distributions before and after learning. Consequently, if we compare the sampling of activity patterns in pre-task sleep to sampling in post-task sleep, then any patterns with changes in their sampling should be from the updated model. This means that these patterns should encode some aspect of the task.

Remarkably, this is exactly what we found. For each co-activation pattern, we found its ability to predict a trial’s outcome by its rate of occurrence on that trial (Fig 5A). When we compared this outcome prediction to the change in sampling between pre- and post-task sleep, we found a strong correlation between the two (Fig. 5B-D). This correlation was highly robust (Fig. 5E-G). The learnt internal model, as evidenced by the updated patterns sampled from it, was seemingly encoding the task.

### Outcome-predictive patterns occur around the choice point.

Consistent with the internal model being task-related, we further found that the outcome-predictive activity patterns preferentially occurred around the choice point of the maze (Fig. 6). Particularly striking was that patterns strongly predictive of outcome rarely occurred in the starting arm (Fig. 6A). Together, the selective changes over learning to outcome-specific (Fig. 5) and location-specific (Fig. 6) activity patterns show that the convergence of distributions (Fig. 2) is not a statistical curiosity, but is evidence for the updating of a behaviourally-relevant internal model.

## Discussion

Prefrontal cortex has been implicated in both planning and working memory during spatial navigation [25–28], and executive control in general [29, 30]. Our results suggests a probabilistic basis for these functions. We find that moment-to-moment patterns of mPFC population activity change their

sampling rates during learning of a spatial navigation task. Consequently, the statistical distributions of patterns in spontaneous and task-evoked activity converge. Our analyses thus suggest mPFC encodes a probabilistic internal model of a task, which is updated by behavioural outcomes, and uses population-activity sampling as the basis for inference.

Remarkably we observed the convergence of distributions using precise activity patterns down to 2 ms resolution. Using surrogate models, we showed that the convergence is due to changes in correlations between neurons, rather than changes in firing rates. Previous work observed fine structure in stimulus-evoked population activity patterns in retina [e.g. 22, 23] and V1 [11]. We extend these results to show that such fine time-scale correlation structure can be observed in cortical regions for executive control, and be evoked by tasks. Unexpectedly, we have shown that, despite their high temporal resolution, task information can be decoded from these patterns.

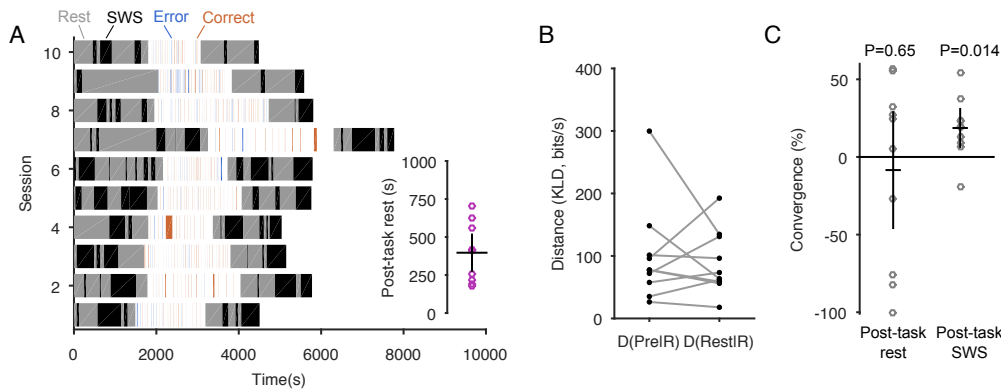
How a cortical region encodes an internal model is an intriguing open question. A strong candidate is the relative strengths of the synaptic connections both into and within the encoding cortical circuit [7, 8, 10, 12]. The activity of a cortical circuit is strongly dependent on the pattern and strength of the connections between its neurons [e.g. 31, 32]. Consequently, defining the underlying model as the circuit’s synaptic network allows both model-based inference through synaptically-driven activity and model learning through synaptic plasticity [10].

Our results are distinct from previous observations of task-specific replay during sleep in prefrontal cortex [33], including reports [16] using the same data analysed here. In contrast to the work here, replay accounts do not consider the statistical distributions of the observed patterns, nor identify the changed patterns [beyond example templates in ref. 33], nor relate them to task behaviour; moreover, replay is described for coincident activity on coarse time-scales greater than those used here by a factor of 50 ([ref. 16]) up to a factor of 10000 ([ref. 33]). They thus do not address the statistical changes to population-wide activity predicted by theories of probabilistic population coding.

Our theory proposes that spontaneous neural activity during sleep is sampling a prior distribution generated from an internal model. We found that the set of activity patterns was remarkably conserved between sleeping and behaviour (Fig. S1), despite the different global dynamics of cortex between these states [34, 35], consistent with activity being generated from the same internal model in both states. This theory predicts that manipulating synaptic weights during sleep, changing the internal model, should change both the prior and the posterior distributions over task variables. Recent work has shown that inducing task-specific reward signals during sleep, likely altering synaptic weights, indeed immediately alters task behaviour on waking [36]. Our results thus suggest that casting sleeping and waking activity as prior and posterior distributions generated from the same internal model could be a fruitful computational framework for relating cortical dynamics to behaviour.

## Materials and Methods

**Task and electrophysiological recordings.** The data analysed here were from ten learning sessions and seven rule change sessions in



**Fig. 4.** Convergence is not a recency effect. (A) Breakdown of each learning session into the duration of its state components. The task epoch is divided into correct (red) and error (blue) trials, and inter-trial intervals (white spaces). Trial durations were typically 2-4 seconds, so are thin lines on this scale. The pre- and post-task epochs contained quiet waking and light sleep states ("Rest" period) and identified bouts of slow-wave sleep ("SWS"). Inset: duration of the Rest period between the end of the last trial and the start of the first SWS bout (lines give mean  $\pm$  2 s.e.m.) (B) Distances between the distributions of pattern frequencies in different epochs; one dot per session.  $D(X|Y)$ : distance between pattern distributions in epochs  $X$  and  $Y$ : Pre: pre-task SWS; Rest: immediate post-task rest period; R: correct task trials. Compare to Fig. 2A. (C) Results from panel B expressed as the convergence between the distributions in the task and post-task rest period. We also re-plot here the convergence between the task and post-task SWS distributions. ( $P$ -values from 1-tailed Wilcoxon signrank test, with  $N=10$  sessions).

the study of [ref. 16]. For full details on training, spike-sorting, and histology see [16]. Four Long-Evans male rats with implanted tetrodes in prelimbic cortex were trained on the Y-maze task (Fig. 1a). Each recording session consisted of a 20-30 minute sleep or rest epoch (pre-task epoch), in which the rat remained undisturbed in a padded flowerpot placed on the central platform of the maze, followed by a task epoch, in which the rat performed for 20-40 minutes, and then by a second 20-30 minute sleep or rest epoch (post-task epoch). Every trial started when the rat reached the departure arm and finished when the rat reached the end of one of the choice arms. Correct choice was rewarded with drops of flavoured milk. Each rat had to learn the current rule by trial-and-error, either: go to the left arm; go to the right arm; go to the lit arm. To maintain consistent context across all sessions, the extra-maze light cues were lit in a pseudo-random sequence across trials, whether they were relevant to the rule or not.

We primarily analysed here data from the ten sessions in which the previously-defined learning criteria were met: the first trial of a block of at least three consecutive rewarded trials after which the performance until the end of the session was above 80%. In later sessions the rats reached the criterion for changing the rule: ten consecutive correct trials or one error out of 12 trials. Thus each rat learnt at least two rules, with seven rule-change sessions in total.

Tetrode recordings were spike-sorted only within each recording session for conservative identification of stable single units. In the 17 sessions we analyse here, the populations ranged in size from 15-55 units.

**Activity pattern distributions.** For a population of size  $N$ , we characterised population activity from time  $t$  to  $t+\delta$  as an  $N$ -length binary vector with each element being 1 if at least one spike was fired by that neuron in that time-bin, and 0 otherwise. In the main text we use a binsize of  $\delta = 2$  ms throughout, and report the results of using larger binsizes in (Fig. S2). We build patterns using the number of recorded neurons  $N$ , up to a maximum of 35 for computational tractability. The probability distribution for these activity patterns was compiled by counting the frequency of each pattern's occurrence and normalising by the total number of pattern occurrences.

**Comparing distributions.** We quantified the distance  $D(P|Q)$  between probability distributions  $P$  and  $Q$  using both the Kullback-Liebler divergence (KLD) and the Hellinger distance.

The KLD is an information theoretic measure to compare the similarity between two probability distributions. Let  $P = (p_1, p_2, \dots, p_n)$  and  $Q = (q_1, q_2, \dots, q_n)$  be two discrete probability distributions, for  $n$  distinct possibilities – for us, these are all possible individual activity patterns. The KLD is then defined as  $d(P|Q) = \sum_{i=1}^n p_i \ln(\frac{p_i}{q_i})$ . This measure is not symmetric, so that in general  $d(P|Q) \neq d(Q|P)$ . Following prior

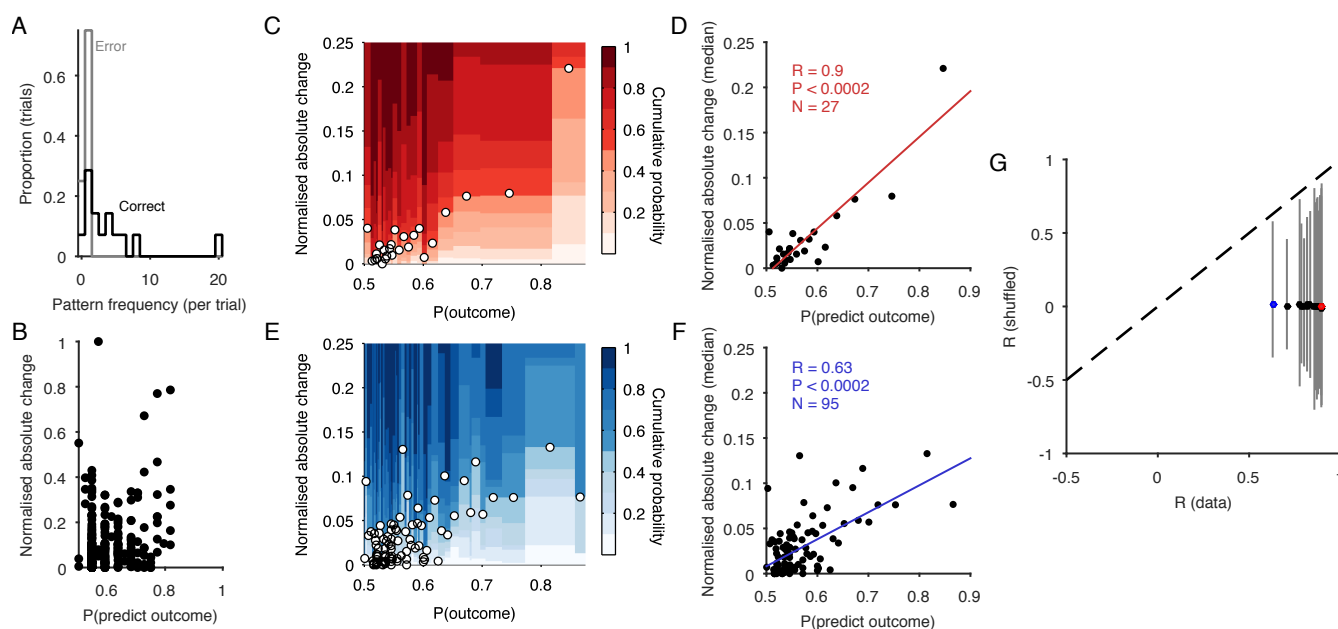
work [11, 20], we thus compute and report the symmetrised KLD:  $D(P|Q) = (d(P|Q) + d(Q|P))/2$ .

There are  $2^N$  distinct possible activity patterns in a recording with  $N$  neurons. Most of these activity patterns are never observed, so we exclude the activity patterns that are not observed in either of the epochs we compare. The empirical frequency of the remaining activity patterns is biased due to the limited length of the recordings [19]. To counteract this bias, we use the Bayesian estimator and quadratic bias correction exactly as described in [11]. The Berkes estimator assumes a Dirichlet prior and multinomial likelihood to calculate the posterior estimate of the KLD; we use their code ([github.com/pberkes/neuro-kl](https://github.com/pberkes/neuro-kl)) to compute the estimator. We then compute a KLD estimate using all  $S$  activity patterns, and using  $S/2$  and  $S/4$  patterns randomly sampled without replacement. By fitting a quadratic polynomial to these three KLD estimates, we can then use the intercept term of the quadratic fit as an estimate of the KLD if we had access to recordings of infinite length [19, 37].

The Hellinger distance for two discrete distributions  $P$  and  $Q$  is  $D(P|Q) = \frac{1}{2} \sum_{i=1}^n (\sqrt{p_i} - \sqrt{q_i})^2$ . To a first approximation, this measures for each pair of probabilities  $(p_i, q_i)$  the distance between their square-roots. In this form,  $D(P|Q) = 0$  means the distributions are identical, and  $D(P|Q) = 1$  means the distributions are mutually singular: all positive probabilities in  $P$  are zero in  $Q$ , and vice-versa. The Hellinger distance is a lower bound for the KLD:  $2D(P|Q) \leq KLD$ .

To compare distances between sessions we computed a normalised measure of "convergence". The divergence between a given pair of distributions could depend on many factors that differ between sessions, including that each recorded population was a different size, and how much of the relevant population for encoding the internal model we recorded. Consequently, the key comparison between the divergences  $D(Pre|R) - D(Post|R)$  also depends on these factors. To compare the difference in divergences across sessions, we computed a "convergence" score by normalising by the scale of the divergence in the pre-task SWS:  $((D(Pre|R) - D(Post|R)) / D(Pre|R))$ . We express this as a percentage. Convergence greater than 0% indicates that the distance between the task ( $R$ : correct trials) and post-task SWS ( $Post$ ) distributions is smaller than that between the task and pre-task SWS ( $Pre$ ) distributions.

**Statistics.** Quoted measurement values are means  $\pm$  s.e.m. All hypothesis tests used the non-parametric Wilcoxon signrank test for a one-sample test that the sample median for the population of sessions is greater than zero. For learning sessions, we have  $N = 10$  sessions; for rule-changes (Fig. 2) we have  $N = 7$  sessions. Throughout we plot mean values and their approximate 95% confidence intervals given by  $\pm 2$  s.e.m.



**Fig. 5.** Coding of trial outcome by sampled activity patterns. (A) Example distributions of a pattern's frequency conditioned on trial outcome from one session. (B) For all co-activation patterns in one session, a scatter plot of outcome prediction and (absolute) change in pattern frequency between pre- and post-task SWS. Change is normalised to the maximum change in the session. (C) Distribution of change in pattern frequency according to outcome prediction over all ten sessions. Colour intensity gives the cumulative probability of at least that change. Circles give the median absolute change for each distribution. In this example, distributions were built using bins with 90 data-points each. Unbinned data are analysed in Fig. S6. (D) Correlation of outcome prediction and median change in pattern occurrence between sleep epochs from C, over all ten sessions. Red line is the best-fit linear regression ( $P < 0.0002$ , permutation test). (E)-(F) As C-D, for the worst-case correlation observed, with 25 data-points per bin. (G) Robustness of correlation results. Solid dots plot the correlation coefficient  $R$  between outcome prediction and median change in pattern frequency obtained for different binnings of the data. Coloured dots correspond to panels C-D and E-F. Lines each give the entire range of  $R$  obtained from a 5000-repeat permutation test; none reach the equivalent data point (dashed line shows equality), indicating all data correlations had  $P < 0.0002$ .

Bootstrapped confidence intervals (in Fig. 2C,H) for each session were constructed using 1000 bootstraps of each epoch's activity pattern distribution. Each bootstrap was a sample-with-replacement of activity patterns from the data distribution  $X$  to get a sample distribution  $X^*$ . For a given pair of bootstrapped distributions  $X^*, Y^*$  we then compute their distance  $D^*(X^*|Y^*)$ . Given both bootstrapped distances  $D^*(Pre|R)$  and  $D^*(Post|R)$ , we then compute the bootstrapped convergence  $(D^*(Pre^*|R^*) - D^*(Post^*|R^*)) / D^*(Pre^*|R^*)$ .

**Raster model.** To control for the possibility that changes in activity pattern occurrence were due solely to changes in the firing rates of individual neurons and the total population, we used the raster model exactly as described in [20]. For a given data-set of spike-trains  $N$  and binsize  $\delta$ , the raster model constructs a synthetic set of spikes such that each synthetic spike-train has the same mean rate as its counterpart in the data, and the distribution of the total number of spikes per time-bin matches the data. In this way, it predicts the frequency of activity patterns that should occur given solely changes in individual and population rates.

For Fig 3 we generated 1000 raster models per session using the spike-trains from the post-task SWS in that session. For each generated raster model, we computed the distance between its distribution of activity patterns and the data distribution for correct trials in the task  $D(Post - model|R)$ . This comparison gives the expected distance between task and post-task SWS distributions due to firing rate changes alone. We plot the difference between the mean  $D(Post - model|R)$  and the data  $D(Post|R)$  in Fig. 3.

**Outcome prediction.** We examined the correlates of activity pattern occurrence with behaviour. To rule out pure firing rate effects, we excluded all patterns with  $K = 0$  and  $K = 1$  spikes, considering only co-activation patterns with two or more active neurons.

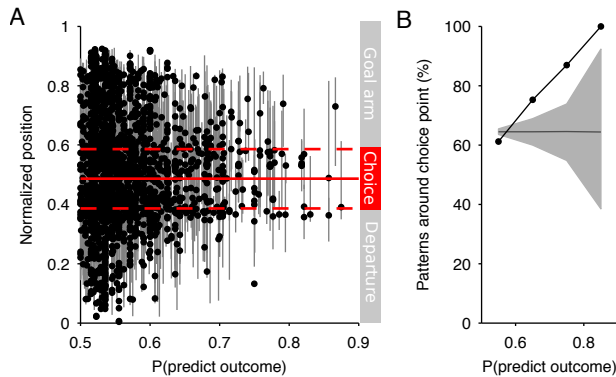
To check whether individual activity patterns coded for the outcome on each trial, we used standard receiver-operating characteristic (ROC) analysis. For each pattern, we computed the distribution of its occurrence frequencies separately for correct and

error trials (as in the example of Fig. 5A). We then used a threshold  $T$  to classify trials as error or correct based on whether the frequency on that trial exceeded the threshold or not. We found the fraction of correctly classified correct trials (true positive rate) and the fraction of error trials incorrectly classified as correct trials (false positive rate). Plotting the false positive rates against the true positive rates for all values of  $T$  gives the ROC curve. The area under the ROC curve gives the probability that a randomly chosen pattern frequency will be correctly classified as from a correct trial; we report this as  $P(\text{predict outcome})$ .

**Relationship of sampling change and outcome prediction.** Within each session, we computed the change in each pattern's occurrence between pre- and post-task SWS. These were normalised by the maximum change within each session. Maximally changing patterns were candidates for those updated by learning during the task. Correlation between change in pattern sampling and outcome prediction was done on normalised changes pooled over all sessions. Change scores were binned using variable-width bins of  $P(\text{predict outcome})$ , each containing the same number of data-points to rule out power issues affecting the correlation. We regress  $P(\text{predict outcome})$  against the median change in each bin, using the mid-point of each bin as the value for  $P(\text{predict outcome})$ . Our main claim is that prediction and change are dependent variables (Fig. 5C-G). To test this claim, we compared the data correlation against the null model of independent variables, by permuting the assignment of change scores to the activity patterns. For each permutation, we repeat the binning and regression. We permuted 5000 times to get the sampling distribution of the correlation coefficient  $R^*$  predicted by the null model of independent variables. To check robustness, all analyses were repeated for a range of fixed number of data-points per bin between 20 and 100.

**Relationship of location and outcome prediction.** The location of every occurrence of a co-activation pattern was expressed as a normalized position on the linearised maze (0: start of departure arm;

1: end of the chosen goal arm). Our main claim is that activity patterns strongly predictive of outcome occur predominantly around the choice point of the maze, and so prediction and overlap of the choice area are dependent variables (Fig. 6B). To test this claim, we compared this relationship against the null model of independent variables, by permuting the assignment of location centre-of-mass (median and interquartile range) to the activity patterns. For each permutation, we compute the proportion of patterns whose interquartile range overlaps the choice area, and bin as per the data. We permuted 5000 times to get the sampling distribution of the proportions predicted by the null model of independent variables:



**Fig. 6.** Outcome predicting activity patterns are sampled in the choice area. (A) Scatter plot of each pattern's outcome prediction and sample locations in the maze (dot is median position; grey line is interquartile range); all positions given as a proportion of the linearised maze from start of departure arm. Red lines indicate the approximate centre (solid) and boundaries (dashed) of the maze's choice area (cf Fig 1A). (B) Proportion of activity patterns whose interquartile range of sample locations enters the choice area (black dots and line). Grey region shows mean (line) and 95% range (shading) of proportions from a permutation test. The data exceed the upper limit of expected proportions for all outcome-predictive patterns.

- Kording KP, Wolpert DM (2004) Bayesian integration in sensorimotor learning. *Nature* 427:244–247.
- Pouget A, Beck JM, Ma WJ, Latham PE (2013) Probabilistic brains: knowns and unknowns. *Nat Neurosci* 16:1170–1178.
- Wolpert DM, Ghahramani Z, Jordan MI (1995) An internal model for sensorimotor integration. *Science* 269:1880–1882.
- Dayan P, Abbot LF (2001) *Theoretical Neuroscience* ed. anonymous. (MIT Press, Cambridge, MA).
- Zemel RS, Dayan P, Pouget A (1998) Probabilistic interpretation of population codes. *Neural Comput* 10:403–430.
- Ma WJ, Beck JM, Latham PE, Pouget A (2006) Bayesian inference with probabilistic population codes. *Nat Neurosci* 9:1432–1438.
- Buesing L, Bill J, Nessler B, Maass W (2011) Neural dynamics as sampling: a model for stochastic computation in recurrent networks of spiking neurons. *PLoS Comput Biol* 7:e1002211.
- Kappel D, Habenschuss S, Legenstein R, Maass W (2015) Network plasticity as Bayesian inference. *PLoS Comput Biol* 11:e1004485.
- Beck JM et al. (2008) Probabilistic population codes for Bayesian decision making. *Neuron* 60:1142–1152.
- Fiser J, Berkes P, Orbán G, Lengyel M (2010) Statistically optimal perception and learning: from behavior to neural representations. *Trends Cogn Sci* 14:119–130.
- Berkes P, Orbán G, Lengyel M, Fiser J (2011) Spontaneous cortical activity reveals hallmarks of an optimal internal model of the environment. *Science* 331:83–87.
- Habenschuss S, Jonke Z, Maass W (2013) Stochastic computations in cortical microcircuit

we plot the mean and 95% range of this sampling distribution as the grey region in Fig. 6B.

**ACKNOWLEDGMENTS.** We thank the Humphries lab (Javier Caballero, Mat Evans, Silvia Maggi) for discussions; Rasmus Petersen for comments on the manuscript; and P. Berkes and M. Okun for respectively making their KL divergence and raster model code publicly available. A.S. and M.D.H were supported by a Medical Research Council Senior non-Clinical Fellowship award to M.D.H. A.P. was supported by Human Frontier Science Program Fellowship LT000160/2011-1 and National Institute of Health Award K99 NS086915-01.

- models. *PLoS Comput Biol* 9:e1003311.
- Ragozzino ME, Detrick S, Kesner RP (1999) Involvement of the prelimbic-infralimbic areas of the rodent prefrontal cortex in behavioral flexibility for place and response learning. *J Neurosci* 19:4585–4594.
- Rich EL, Shapiro ML (2007) Prelimbic/infralimbic inactivation impairs memory for multiple task switches, but not flexible selection of familiar tasks. *J Neurosci* 27:4747–4755.
- Benchenane K et al. (2010) Coherent theta oscillations and reorganization of spike timing in the hippocampal- prefrontal network upon learning. *Neuron* 66:921–936.
- Peyrache A, Khamassi M, Benchenane K, Wiener SI, Battaglia FP (2009) Replay of rule-learning related neural patterns in the prefrontal cortex during sleep. *Nat Neurosci* 12:916–926.
- Luczak A, Barthó P, Harris KD (2009) Spontaneous events outline the realm of possible sensory responses in neocortical populations. *Neuron* 62:413–425.
- Wohrer A, Humphries MD, Machens C (2013) Population-wide distributions of neural activity during perceptual decision-making. *Prog Neurobiol* 103:156–193.
- Panzeri S, Senatore R, Montemurro MA, Petersen RS (2007) Correcting for the sampling bias problem in spike train information measures. *J Neurophysiol* 98:1064–1072.
- Okun M et al. (2012) Population rate dynamics and multineuron firing patterns in sensory cortex. *J Neurosci* 32:17108–17119.
- Fiser J, Lengyel M, Savin C, Orbán G, Berkes P (2013) How (not) to assess the importance of correlations for the matching of spontaneous and evoked activity. p. arXiv:1301.6554.
- Schneidman E, Berry MJ, Segev R, Bialek W (2006) Weak pairwise correlations imply strongly correlated network states in a neural population. *Nature* 440:1007–1012.
- Tkacik G et al. (2014) Searching for collective behavior in a large network of sensory neurons. *PLoS Comput Biol* 10:e1003408.
- Ganmor E, Segev R, Schneidman E (2015) A thesaurus for a neural population code. *Elife* 4:e06134.
- Baeg EH et al. (2003) Dynamics of population code for working memory in the prefrontal cortex. *Neuron* 40:177–188.
- Fujisawa S, Amarasingham A, Harrison MT, Buzsáki G (2008) Behavior-dependent short-term assembly dynamics in the medial prefrontal cortex. *Nat Neurosci* 11:823–833.
- Ito HT, Zhang SJ, Witter MP, Moser EI, Moser MB (2015) A prefrontal-thalamo-hippocampal circuit for goal-directed spatial navigation. *Nature* 522:50–55.
- Spellman T et al. (2015) Hippocampal-prefrontal input supports spatial encoding in working memory. *Nature* 522:309–314.
- Miller EK (2000) The prefrontal cortex and cognitive control. *Nat Rev Neurosci* 1(1):59–65.
- Sul JH, Kim H, Huh N, Lee D, Jung MW (2010) Distinct roles of rodent orbitofrontal and medial prefrontal cortex in decision making. *Neuron* 66:449–460.
- Cossell L et al. (2015) Functional organization of excitatory synaptic strength in primary visual cortex. *Nature* 518:399–403.
- Okun M et al. (2015) Diverse coupling of neurons to populations in sensory cortex. *Nature* 521:511–515.
- Euston DR, Tatsuno M, McNaughton BL (2007) Fast-forward playback of recent memory sequences in prefrontal cortex during sleep. *Science* 318:1147–1150.
- Destexhe A, Contreras D, Steriade M (1999) Spatiotemporal analysis of local field potentials and unit discharges in cat cerebral cortex during natural wake and sleep states. *J Neurosci* 19:4595–4608.
- Steriade M, Timofeev I, Grenier F (2001) Natural waking and sleep states: a view from inside neocortical neurons. *J Neurophysiol* 85:1969–1985.
- de Lavilleon G, Lacroix MM, Rondi-Reig L, Benchenane K (2015) Explicit memory creation during sleep demonstrates a causal role of place cells in navigation. *Nat Neurosci* 18:493–495.
- Strong SP, Koberle R, de Ruyter van Steveninck RR, Bialek W (1998) Entropy and information in neural spike trains. *Phys Rev Lett* 80:197–200.



## SI Text

### Neural inference

How do we know the current state of the world given some input from it? Our input is both limited in time and noisy, so our estimates are inherently uncertain. Consequently, we have an inference problem: what is our best guess of the current state of the world given some finite, noisy input? We can state this problem as being equivalent to inferring the probability distribution

$$P(\text{state}|\text{input}, \text{model}) \quad [1]$$

at some given moment in time  $t$ ; in words, this is the probability of currently being in a given state, out of all possible states, given both the available input and some internal model of the world. Using Bayes' theorem, we can make this dependence on input and model explicit:

$$\underbrace{P(\text{state}|\text{input}, \text{model})}_{\text{posterior}} \propto \underbrace{P(\text{input}|\text{state}, \text{model})}_{\text{likelihood}} \underbrace{P(\text{state}|\text{model})}_{\text{prior}} \quad [2]$$

The prior is the internal estimate of the current state before the observation *input*, the posterior is the estimate of the current state after observing *input*, and the improvement in the estimate arises from the new information available in *input* that is processed through the likelihood. All these are dependent on the model of the world we are using. This internal model specifies how we interpret the inputs in the likelihood, and generate the prior probabilities. If we change the model, we change these two operations, and so change our estimates of the current state of the world. We can think of the model as specifying what we expect to be relevant in the input, and what states we expect to be in.

One goal of learning is thus to update the internal model to match the statistical properties of the world. The better the model, the better we will be able to predict the state of the external world. But as we can only access directly the inputs generated from those states, formally we say that learning seeks to maximise  $P(\text{input}|\text{model})$  over all possible inputs at all times  $t$  by changing the parameters of the model. A model which always generates maximum values for  $P(\text{input}|\text{model})$  is the best possible learnt internal model of the external world. Obtaining such a model necessarily means that we have experienced all possible states giving rise to those inputs, so that the prior  $P(\text{state}|\text{model})$  is always accurate, and we obtain no new information from the likelihood. Consequently, the posterior probability becomes always proportional to the prior probability. A measure of learning is thus how close the prior and posterior distributions have become.

### Inference-by-sampling

The inference-by-sampling theory [1, 2] proposes that the model is encoded by the particular set and weight of connections in a neural circuit. In this view, the posterior distribution is encoded by the activity of the circuit evoked by some input. Crucially, it predicts that the prior distribution is encoded by spontaneous activity of the same circuit, as this is solely sampling the model.

If the circuit is the model, then the theory predicts that the circuit's instantaneous population activity is a sample from a probability distribution - from the posterior when receiving

external input, from the prior in spontaneous activity. Some downstream neurons, receiving these samples as a consecutive sequence of inputs, can reconstruct the probability distribution just by summing their inputs over time.

For simplicity, Berkes et al [2] considered the instantaneous population activity as some binary vector indicating whether each neuron was active or inactive in a very small time window. This representation makes the distributions easy to measure experimentally.

Learning updates synaptic weights, altering the encoded model. The prediction that posterior and prior distributions converge over learning is thus neurally equivalent to the convergence between the distributions of evoked and spontaneous population activity.

### Evidence for inference-by-sampling in visual cortices

These ideas were developed in the context of visual processing, and particularly with reference to V1. In this context, the "state" of the world is the current view, and the input is the information received by the retina. The proposed purpose of inference in V1 is to infer the most likely low level visual features - edges, for example - present in the current view, given the input to the retina. V1's internal model is then a statistical model of the low-level features, which can be built over a life-time's experience of the world.

Consequently, Berkes and colleagues [2] tested the construction of this internal model by recording from area V1 at different stages of development in the ferret. Natural images were used to probe the current posterior distribution supported by the model, and darkness was used to probe the current prior distribution. Over development, the activity distribution evoked by natural images increased its similarity to the distribution during darkness. This increase was robust to a series of controls for simultaneous changes in firing rate statistics [2-4]. Their results are consistent with the inference-by-sampling interpretation in which the internal model is updated by experience with the world, so that the posterior and prior distributions converge.

### Inference-by-sampling in higher cortices over learning during behaviour

These results could not address learning separately from development. Further, unknown is whether inference-by-sampling can be observed in higher-order cortices, or during ongoing behaviour.

There is no *a priori* reason to expect that inference-by-sampling would be restricted to primary sensory cortices. Much has been written about the generic nature of the cortical microcircuit [5, 6], so we might reasonably expect that, if an internal model is encoded by the neural circuit in V1, so other similar cortical circuits in other regions encode other internal models.

Compelling support for this has come from modelling work by Maass and colleagues [7, 8]. Their models have shown how a wide range of plausible cortical circuit models all produce the necessary dynamics to sample from a statistical model encoded by the circuit's connections [7, 8]. Moreover, the models also replicate key properties of the firing statistics in cortex, including the close-to-Poisson irregularity of firing patterns. These suggest that the inference-by-sampling hypothesis is indeed a plausible generic computation for cortex.

Inference of state is also a generic operation. Nothing in Equation 2 limits its application to sensory information. We might consider “state” in the sense used in the reinforcement learning literature [9], as a generic description of the current values of variables of the external world. Indeed, in forms of reinforcement learning that depend on simulation of future actions, “state” in this context can even refer to the simulated values of variables in the external world - for which we would use the internal model to simulate possible outcomes. During behaviour, we might thus expect that an internal model is learnt about the statistical dependence of outcomes on decisions in particular contexts.

The power of the inference-by-sampling hypothesis is that we do not need to know the internal model to test for its existence. We need not specify an exact model to test the convergence of distributions in evoked and spontaneous activity, but such a convergence is evidence of an updated internal model.

Consequently, to test the generality of the inference-by-sampling hypothesis, we sought to test the convergence of distributions over learning using data from the medial prefrontal cortex (mPFC) of rats learning rules in a Y-maze task [10]. By looking at these data for a change to some internal model in mPFC, we are assuming only that the model is related to the rule, but not any specific form of model. It could encode the set of task states and their transitions; it could encode the current sequence of required actions; it could be a statistical model of outcomes. Supporting this assumption, we know mPFC is necessary for successful acquisition of new rules [11, 12], and that mPFC pyramidal neurons change their firing patterns during acquisition of the rules used here [13].

Even if the interpretation of the convergence of distributions in the inference-by-sampling framework turns out to be incorrect, the observation of such a convergence between waking and spontaneous activity over learning still offers compelling clues to the nature of cortical computation.

## What distributions to compare?

Nonetheless, the inference-by-sampling theory places limits on exactly which activity distributions to compare. In the Berkes et al. [2] study, this decision was made simple by the elegant experimental design. As they monitored V1 over development, so it was reasonable to expect the internal model to adapt to the statistics of the world over a lifetime. Their tests at different developmental stages were samples of the current posterior and prior distributions supported by the model. We would not expect significant changes to the internal model during their testing, as it was short on the time-scale of the developmental changes, and so they could compare their entire recorded distributions of evoked and spontaneous activity. In other words, they were able to compare two distributions from the same, static model.

Our data on rats learning rules in a Y-maze allow us to address if learning of the internal model can be observed. But learning on short time-scales brings the confounding issue that learning the model is happening online, while we are monitoring activity. So what distributions should we compare?

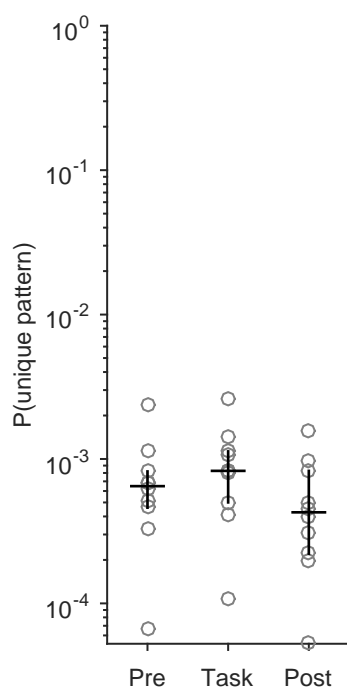
We chose the 10 training sessions in which the rat clearly acquired the present rule, so we could be reasonably sure

that we would observe changes that correlated with learning. We reasoned that neural activity in clearly identified sleep periods before and after the session was a clear candidate for spontaneous activity, as it occurred in the absence of external sensory input. We used slow-wave sleep periods to clearly delineate the presence of sleep. As the rats acquired the rule in that session then, if mPFC indeed encodes rule acquisition, we expect that the spontaneous activity in sleep after the session is drawn from the internal model related to the correct rule.

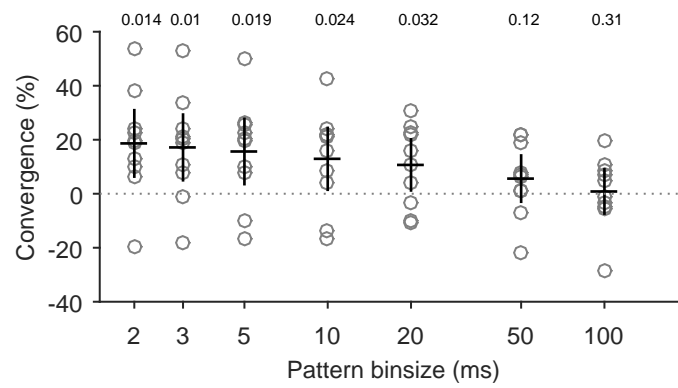
We can only be sure that during behaviour this correct-rule model would be sampled on correct trials. This does not imply that mPFC activity is causal for decisions on those trials - even in a monitoring or goal-encoding role, mPFC activity would reflect whether or not the correct decision was made. The mPFC activity on error trials is unconstrained by the theory. Consequently, we can only be sure that, if the inference-by-sampling hypothesis is true, then the distribution of samples on correct trials would converge, on average, to the distribution in sleep after learning.

The final, subtle constraint is that overt behavioural signs of learning likely indicates ongoing synaptic plasticity. For example, on the same Y-maze, some pyramidal neurons in mPFC change the timing of their spikes in relation to the hippocampal theta rhythm, indicating local circuit plasticity [13]. If so, then the internal model is changing during behaviour. But the internal model putatively sampled in the post-session sleep will be stable. To thus minimise the confound of these changes during behaviour, and compare static posterior and prior distributions [as per 2], we sought to identify where the internal model updating may have finished. A useful proxy for this is the asymptotic behavioural performance. We thus used the trial at which the rat reached the learning criteria as the indicator of relative stability in the internal model. All correct trials from this trial onwards were then used to construct the activity distribution during the task - we call this distribution  $P(R)$  in the main text, and distances measured between it and some other distribution  $P(X)$  we call  $D(X|R)$ .

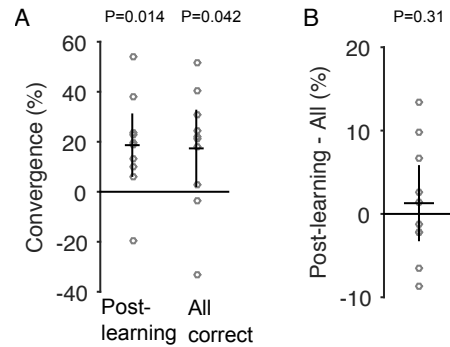
1. Fiser J, Berkes P, Orbán G, Lengyel M (2010) Statistically optimal perception and learning: from behavior to neural representations. *Trends Cogn Sci* 14:119–130.
2. Berkes P, Orbán G, Lengyel M, Fiser J (2011) Spontaneous cortical activity reveals hallmarks of an optimal internal model of the environment. *Science* 331:83–87.
3. Okun M et al. (2012) Population rate dynamics and multineuron firing patterns in sensory cortex. *J Neurosci* 32:17108–17119.
4. Fiser J, Lengyel M, Savin C, Orbán G, Berkes P (2013) How (not) to assess the importance of correlations for the matching of spontaneous and evoked activity. p. arXiv:1301.6554.
5. Thomson AM, Lamy C (2007) Functional maps of neocortical local circuitry. *Front Neurosci* 1:19–42.
6. Harris KD, Shepherd GMG (2015) The neocortical circuit: themes and variations. *Nat Neurosci* 18:170–181.
7. Buesing L, Bill J, Nessler B, Maass W (2011) Neural dynamics as sampling: a model for stochastic computation in recurrent networks of spiking neurons. *PLoS Comput Biol* 7:e1002211.
8. Habenschuss S, Jonke Z, Maass W (2013) Stochastic computations in cortical microcircuit models. *PLoS Comput Biol* 9:e1003311.
9. Sutton RS, Barto AG (1998) *Reinforcement Learning: An Introduction*. (MIT Press, Cambridge, MA).
10. Peyrache A, Khamassi M, Benchenane K, Wiener SI, Battaglia FP (2009) Replay of rule-learning related neural patterns in the prefrontal cortex during sleep. *Nat Neurosci* 12:916–926.
11. Ragozzino ME, Detrick S, Kesner RP (1999) Involvement of the prelimbic-infralimbic areas of the rodent prefrontal cortex in behavioral flexibility for place and response learning. *J Neurosci* 19:4585–4594.
12. Rich EL, Shapiro ML (2007) Prelimbic/infralimbic inactivation impairs memory for multiple task switches, but not flexible selection of familiar tasks. *J Neurosci* 27:4747–4755.
13. Benchenane K et al. (2010) Coherent theta oscillations and reorganization of spike timing in the hippocampal-prefrontal network upon learning. *Neuron* 66:921–936.



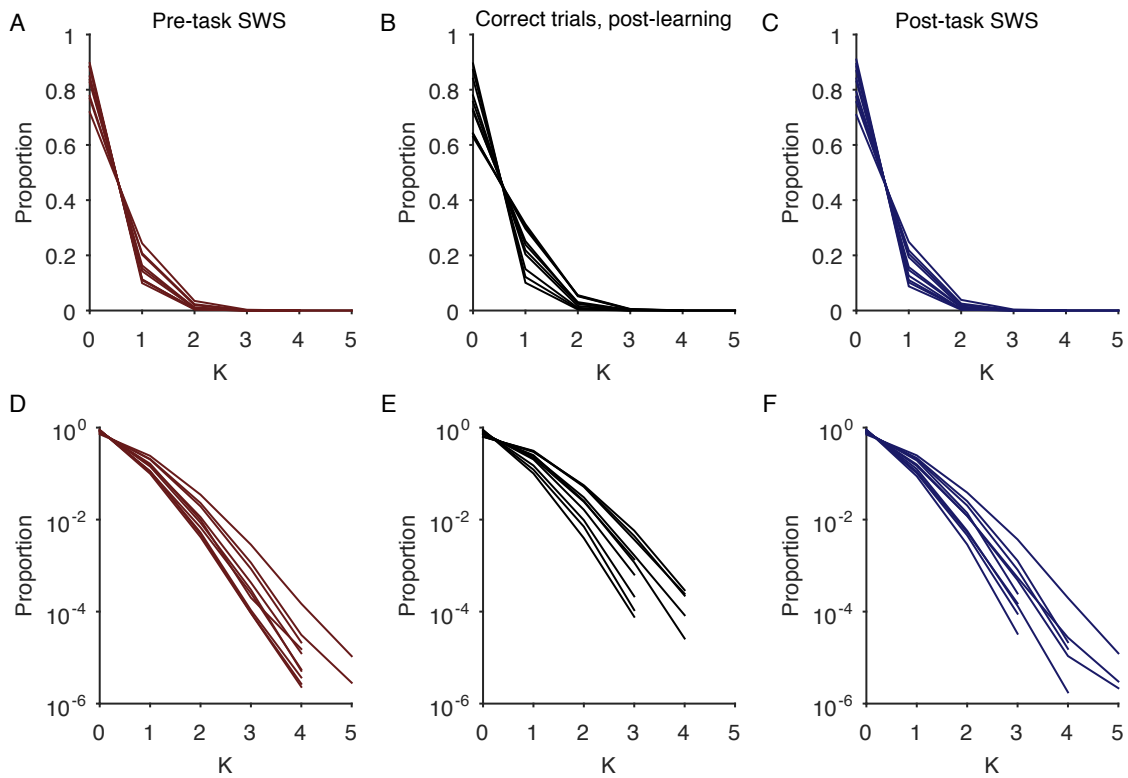
**Fig. S1.** Consistent sampling of activity patterns across session epochs. Each circle is the proportion of activity patterns that appeared only in that epoch of the session. Black bar and line give the median and interquartile range across the 10 sessions. Note the log-scale, showing that the median proportion of unique patterns was less than 0.001 in all three epochs of the session.



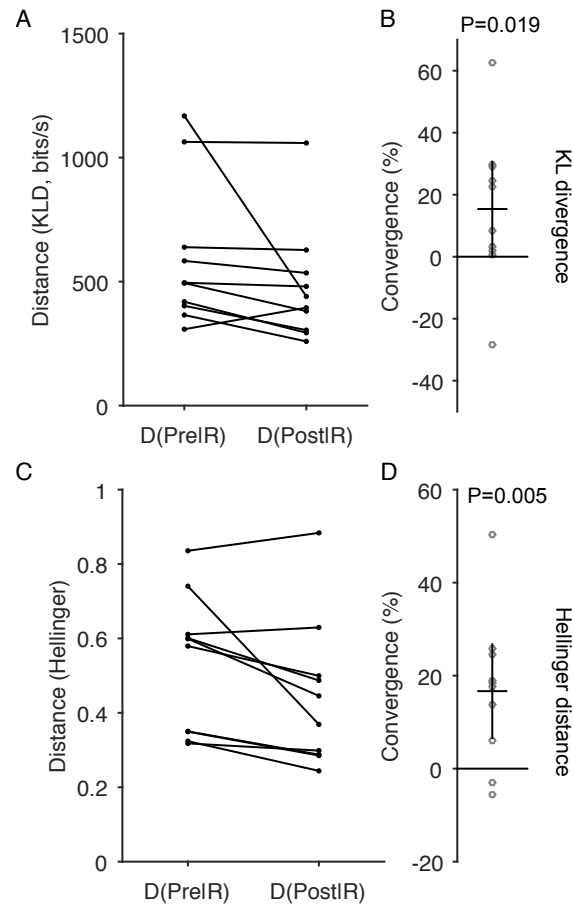
**Fig. S2.** Convergence of distributions over choice of pattern binsize. We plot here the dependence of the convergence of task and post-task SWS distributions on the binsize used for constructing the activity patterns. We see that the convergence of the distribution of patterns on correct task trials is robust to an order of magnitude increase in binsize (the distribution at the binsize of 2ms is plotted in Fig. 2B). Above a binsize of 50 ms, convergence is statistically indistinguishable from zero, meaning that the pre- and post-task SWS distributions are equidistant, on average, from the task distribution. This suggests there is statistical structure in fine time-scale activity patterns that is not present on larger time-scales. Circles are convergence in individual sessions using Kullback-Liebler divergence; black lines give mean  $\pm$  2 s.e.m. Above each distribution is the P-value from a 1-tailed Wilcoxon signrank test, with  $N=10$  sessions.



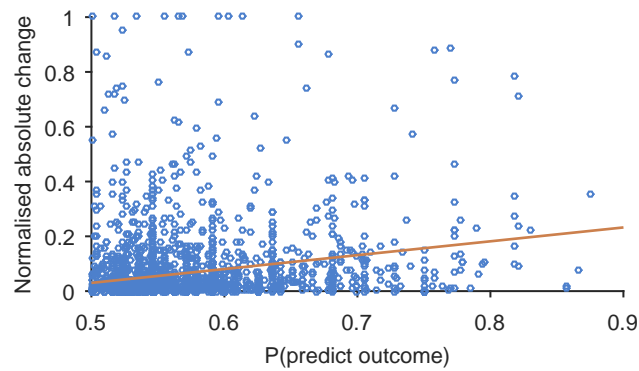
**Fig. S3.** Effect on convergence of including all correct trials. In the main text, we construct all task-related distributions by considering only correct trials after the learning criterion trial (see SI Text). We examine here whether our results were strongly contingent on that choice. (A) If we include all correct trials of a session, we find that convergence between task and post-task SWS distributions is still present ( $P$ -values from 1-tailed Wilcoxon signrank test, with  $N=10$  sessions). Note that the convergence between task and post-task SWS distributions is if anything greater for post-learning trials, even though that task distribution is built from fewer samples, and so might be expected to be noisier. Black lines give mean  $\pm$  2 s.e.m. (B) Difference in convergence between using only post-learning or all correct trials for each session. ( $P$ -values for from 1-tailed Wilcoxon paired-sample ranksum test, with  $N=10$  sessions). Black lines give mean  $\pm$  2 s.e.m.



**Fig. S4.** Distributions of synchronous spiking in all activity patterns. (A)-(C) Distributions of the number of unique recorded activity patterns containing exactly  $K$  spikes, for pre-task SWS (A), correct task trials (B), and post-task SWS (C). Each line is the distribution for one session. (D)-(F) As A-C, plotted on a log-scale to visualise the tails of the distributions. Co-activation patterns ( $K \geq 2$  synchronous spikes) form a small proportion of all patterns.



**Fig. S5.** Convergence between distributions of co-activation patterns. Analysis of distributions restricted to patterns with two or more co-active neurons. (A) Distances between the distributions of pattern frequencies in sleep and task epochs; one dot per session.  $D(X|Y)$ : distance between pattern distributions in epochs  $X$  and  $Y$ : Pre: pre-task SWS; Post: post-task SWS; R: correct task trials. (B) Scatter of convergence across all sessions (circles). Convergence greater than zero means that the activity pattern distribution in the task is closer to the distribution in post-task SWS than the distribution in pre-task SWS. Black lines give mean  $\pm 2$  s.e.m. (C)-(D) As A-B, using Hellinger distance. All  $P$ -values from 1-tailed Wilcoxon signrank test, with  $N=10$  sessions.



**Fig. S6.** Joint distribution of outcome prediction and change in sampling. Here we plot every co-activation pattern's joint values of  $P(\text{outcome})$  and the absolute normalised change in sampling between pre- and post-task slow-wave sleep ( $N = 2353$  patterns with  $K \geq 2$  spikes per pattern across all 10 sessions). The linear regression in red indicates a clear relationship between the two ( $R = 0.22$ ,  $P < 10^{-27}$ ). Nonetheless, the majority of patterns do not markedly change their sampling, nor are they predictive of outcome: 72% (1699/2353) have  $P(\text{outcome}) \leq 0.6$  and a change of less than 10%. Thus fitting a linear regression is not robust, as it is dominated by fitting to this majority that do not change. Rather, it is clear that there is a distribution of change for each  $P(\text{outcome})$ , which we analyse in the main text.