

Prefrontal cortex as a meta-reinforcement learning system

Jane X. Wang^{1,5}, Zeb Kurth-Nelson^{1,2,5}, Dharshan Kumaran^{1,3}, Dhruva Tirumala¹, Hubert Soyer¹, Joel Z. Leibo¹, Demis Hassabis^{1,4} and Matthew Botvinick^{1,4*}

Over the past 20 years, neuroscience research on reward-based learning has converged on a canonical model, under which the neurotransmitter dopamine ‘stamps in’ associations between situations, actions and rewards by modulating the strength of synaptic connections between neurons. However, a growing number of recent findings have placed this standard model under strain. We now draw on recent advances in artificial intelligence to introduce a new theory of reward-based learning. Here, the dopamine system trains another part of the brain, the prefrontal cortex, to operate as its own free-standing learning system. This new perspective accommodates the findings that motivated the standard model, but also deals gracefully with a wider range of observations, providing a fresh foundation for future research.

Exhilarating advances have recently been made toward understanding the mechanisms involved in reward-driven learning. This progress has been enabled in part by the importation of ideas from the field of reinforcement learning¹ (RL). Most centrally, this input has led to an RL-based theory of dopaminergic function. Here, phasic dopamine (DA) release is interpreted as conveying a reward prediction error (RPE) signal², an index of surprise that figures centrally in temporal-difference RL algorithms¹. Under the theory, the RPE drives synaptic plasticity in the striatum, translating experienced action–reward associations into optimized behavioral policies³. Over the past two decades, evidence has steadily mounted for this proposal, establishing it as the standard model of reward-driven learning.

However, even as this standard model has solidified, a collection of problematic observations has accumulated. One quandary arises from research on prefrontal cortex (PFC). A growing body of evidence suggests that PFC implements mechanisms for reward-based learning, performing computations that strikingly resemble those ascribed to DA-based RL. It has long been established that sectors of the PFC represent the expected values of actions, objects and states^{4–6}. More recently, it has emerged that PFC also encodes the recent history of actions and rewards^{5,7–10}. The set of variables encoded, along with observations concerning the temporal profile of neural activation in the PFC, has led to the conclusion that “PFC neurons dynamically [encode] conversions from reward and choice history to object value, and from object value to object choice”⁷. In short, neural activity in PFC appears to reflect a set of operations that together constitute a self-contained RL algorithm.

Placing PFC beside DA, we obtain a picture containing two full-fledged RL systems, one using activity-based representations and the other synaptic learning. What is the relationship between these systems? If both support RL, are their functions simply redundant? One suggestion has been that DA and PFC subservise different forms of learning, with DA implementing model-free RL, based on direct stimulus–response associations, and PFC performing model-based RL, which leverages internal representations of task structure^{11,12}. However, an apparent problem for this dual-system view is the

repeated observation that DA prediction-error signals are informed by task structure, reflecting ‘inferred’^{12,13} and ‘model-based’^{14,15} value estimates that are difficult to square with the standard theory as originally framed.

In the present work we offer a new perspective on the computations underlying reward-based learning, one that accommodates the findings that motivated the existing theories reviewed above, but that also resolves many of the prevailing quandaries. On a wider level, the theory suggests a coherent explanation for a diverse range of findings previously considered unconnected. We begin with three key premises:

System architecture. In line with previous work^{16–18}, we conceptualize the PFC, together with the basal ganglia and thalamic nuclei with which it connects, as forming a recurrent neural network. This network’s inputs include perceptual data, which either contains or is accompanied by information about executed actions and received rewards¹⁹. On the output side, the network triggers actions and also emits estimates of state value (Fig. 1a,b).

Learning. As suggested in past research^{20–22}, we assume that the synaptic weights in the prefrontal network, including its striatal components, are adjusted by a model-free RL procedure, in which DA conveys a RPE signal. Via this role, the DA-based RL procedure shapes the activation dynamics of the recurrent prefrontal network.

Task environment. Following past proposals^{23–25}, we assume that RL takes place not on a single task, but instead in a dynamic environment posing a series of interrelated tasks. The learning system is thus required to engage in ongoing inference and behavioral adjustment.

As indicated, these premises are all firmly grounded in existing research. The new contribution of the present work is to identify an emergent effect that results when the three premises are concurrently satisfied. As we will show, these conditions, when they co-occur, are sufficient to produce a form of ‘meta-learning’²⁶, whereby one learning algorithm gives rise to a second, more efficient learning algorithm. Specifically, by adjusting the connection weights in the prefrontal network, DA-based RL creates a second RL algorithm, implemented entirely in the prefrontal network’s activation dynamics. This new learning algorithm is independent

¹DeepMind, London, UK. ²Max Planck-UCL Centre for Computational Psychiatry and Ageing Research, University College London, London, UK. ³Institute of Cognitive Neuroscience, University College London, London, UK. ⁴Gatsby Computational Neuroscience Unit, University College London, London, UK.

⁵These authors contributed equally: Jane X. Wang and Zeb Kurth-Nelson. *e-mail: botvinick@google.com

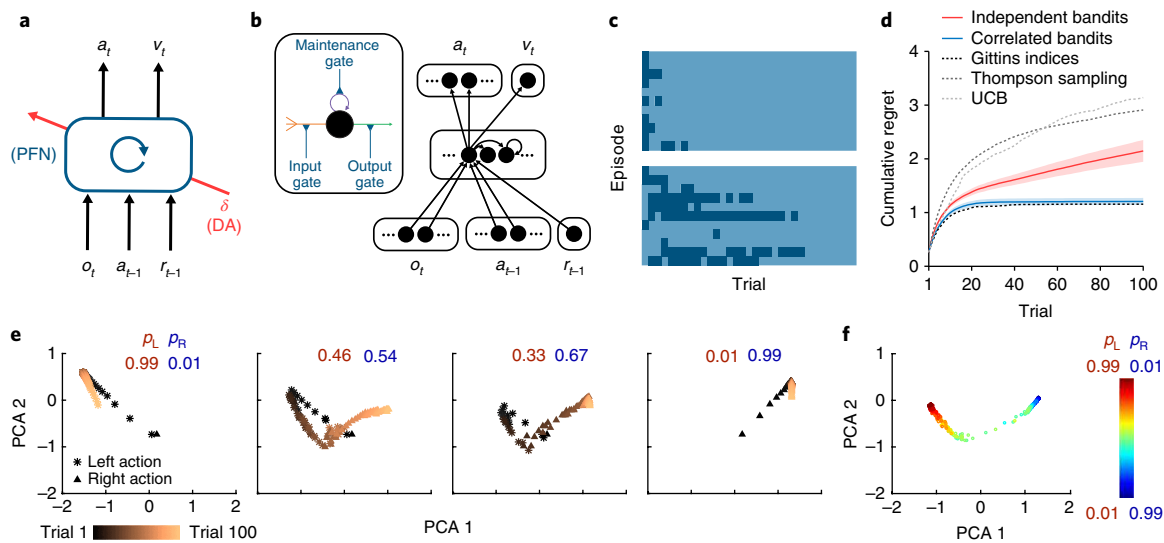


Fig. 1 | Meta-RL architecture learns across episodes to learn efficiently within an episode. **a**, Agent architecture. The prefrontal network (PFN), including sectors of the basal ganglia and the thalamus that connect directly with PFC, is modeled as a recurrent neural network, with synaptic weights adjusted through an RL algorithm driven by DA; o is perceptual input, a is action, r is reward, v is state value, t is time-step and δ is RPE. The central box denotes a single, fully connected set of LSTM units. **b**, A more detailed schematic of the neural network implementation used in our simulations. Input units encoding the current observation and previous action and reward connect all-to-all with hidden units, which are themselves fully connected LSTM units. These units connect all-to-all, in turn, with a softmax layer (see Methods) encoding actions and a single linear unit encoding estimated state value. Inset: a single LSTM unit (see Methods). Input (orange) is a weighted sum of other unit outputs, plus the activity of the LSTM unit itself (purple). Output (green) puts these summed inputs through a sigmoid nonlinearity. All three quantities are multiplicatively gated (blue). Full details, including relevant equations, are presented in Methods. **c**, Trial-by-trial model behavior on bandit problems with Bernoulli (arm 1, arm 2) reward parameters 0.25, 0.75 (top) and 0.6, 0.4 (bottom). The two colors indicate left and right actions, respectively. The network shifts from exploration to exploitation, making this transition more slowly in the more difficult problem. **d**, Performance for the meta-RL network trained on bandits with independently, identically distributed arm parameters and tested on 0.25, 0.75 (red), measured in terms of cumulative regret, defined as the cumulative loss (in expected rewards) suffered when playing suboptimal (lower-reward) arms. Performance of several standard machine-learning bandit algorithms is plotted for comparison. Blue, performance on the same problem after training on correlated bandits (parameters always summed to 1). Shading represents 95% confidence interval over 300 evaluation episodes. **e**, Evolution of recurrent neural network activation pattern during individual trials while testing on correlated bandits after training on problems from the same distribution. p_L , probability of reward for action ‘left’; p_R = probability of reward for action ‘right’. **f**, Recurrent neural network activity patterns from step 100 in the correlated bandit task across a range of payoff parameters. Further analyses are presented in Supplementary Fig. 2. Analysis in **e,f** done for 300 evaluation episodes, each consisting of 100 trials, performed by 1 fully trained network.

of the original one and differs in ways that are suited to the task environment. Crucially, the emergent algorithm is a full-fledged RL procedure: it copes with the exploration–exploitation tradeoff^{27,28}, maintains a representation of the value function¹, and progressively adjusts the action policy. In view of this point, and in recognition of some precursor research^{29–32}, we refer to the overall effect as ‘meta-reinforcement learning’.

For demonstration, we leverage the simple model shown in Fig. 1a, a recurrent neural network (Fig. 1b) whose weights are trained using a model-free RL algorithm, exploiting recent advances in deep learning research (see Methods and Supplementary Fig. 1). We consider this model’s performance in a simple ‘two-armed bandit’ RL task³¹. On each trial, the system outputs an action: left or right. Each has a probability of yielding a reward, but these probabilities change with each training episode, thus presenting a new bandit problem. After training on a series of problems, the weights in the recurrent network are fixed and the system is tested on further problems. The network explores both arms, gradually homing in on the richer one, learning with an efficiency that rivals standard machine-learning algorithms (Fig. 1c,d).

Because the weights in the network were fixed at test, the system’s learning ability cannot be attributed to the RL algorithm that was used to tune the weights. Instead, learning reflects the activation dynamics of the recurrent network. As a result of training, these dynamics implement their own RL algorithm, integrating reward information over time, exploring, and refining the

action policy (Fig. 1e,f and Supplementary Fig. 2). This learned RL algorithm not only functions independently of the algorithm that was originally used to set the network weights; it also differs from that original algorithm in ways that make it specially adapted to the task distribution on which the system was trained. An illustration of this point is presented in Fig. 1d, which shows performance of the same system after training on a structured version of the bandit problem in which the arm parameters were anticorrelated across episodes. Here the recurrent network converges on an RL algorithm that exploits the problem’s structure, identifying the superior arm more rapidly than in the unstructured task.

Having introduced meta-RL in abstract computational terms, we now return to its neurobiological interpretation. This starts by regarding the prefrontal network, including its subcortical components, as a recurrent neural network. DA, as in the standard model, broadcasts an RPE signal, driving synaptic learning in the prefrontal network. The principal role of this learning is to shape the dynamics of the prefrontal network by tuning its recurrent connectivity. Through meta-RL, these dynamics come to implement a second RL algorithm, which differs from the original DA-driven algorithm, assuming a form tailored to the task environment. The role of DA-driven RL, under this account, plays out across extended series of tasks. Rapid within-task learning is mediated primarily by the emergent RL algorithm inherent in the dynamics of the prefrontal network.

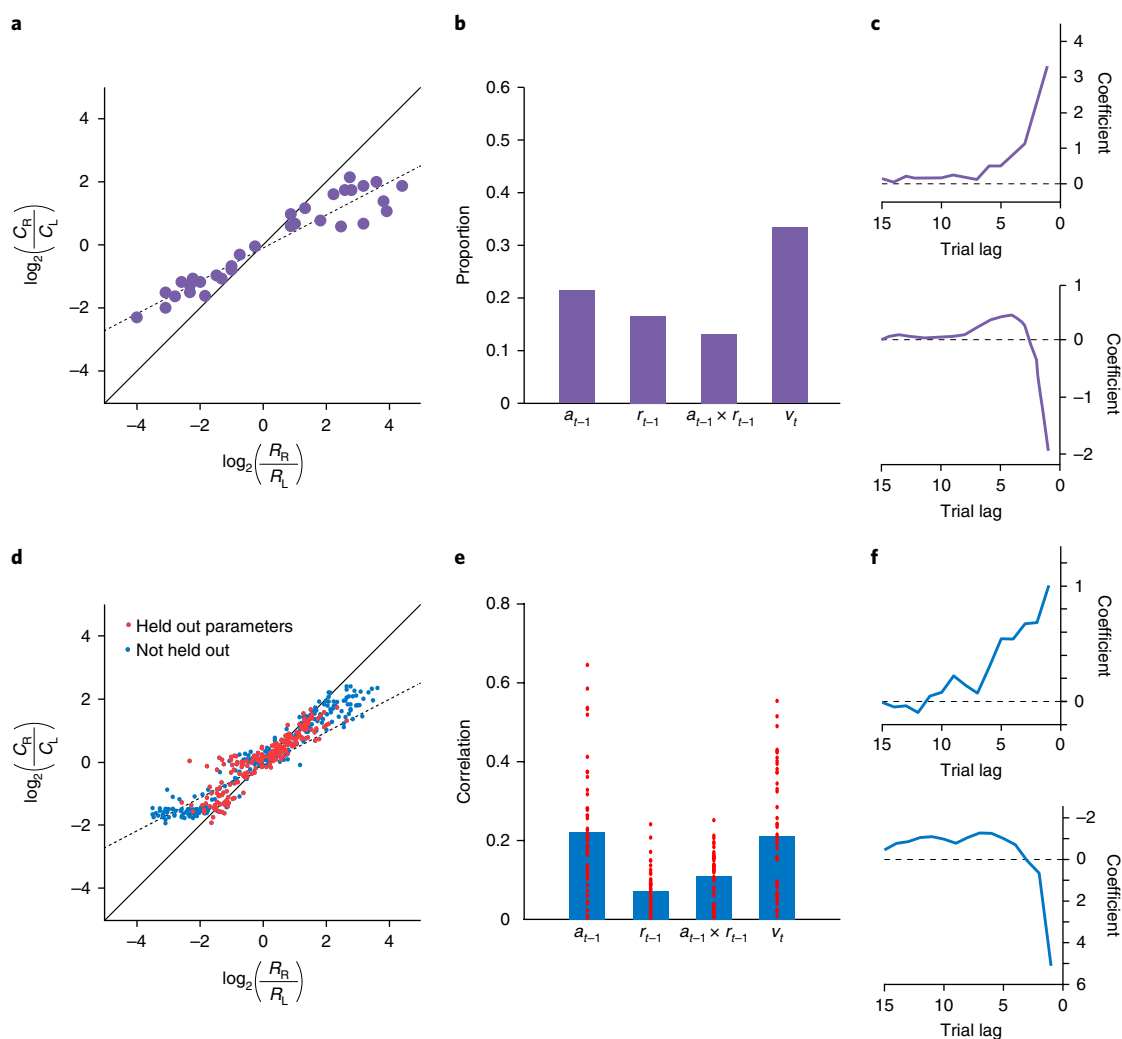


Fig. 2 | Individual units in simulated recurrent network code for action and reward history. **a**, Probability matching behavior from Lau and Glimcher³³. C_i , number of trials on which action i was selected; R_i , number of reward events yielded by action i . **b**, Proportion of PFC neurons encoding, during the trial's initial fixation period, preceding action, preceding reward, the interaction of these two factors, and current choice value, from Tsutsui et al.⁷ **c**, Lag regression coefficients indicating the influence of recent reward outcomes (top) and actions (bottom) on choice from Lau and Glimcher³³. **d**, Matching behavior from model. Red points show generalization performance (see Methods). **e**, Partial correlation coefficients (absolute values, Pearson's r) for the same factors as in **b**, but across units in the model at the trial's initial fixation. Bars indicate mean values. Coding for the upcoming action (not shown) was also ubiquitous across units, and generally stronger than for the other variables. **f**, Regression weights corresponding to those in **c**, but based on model behavior. Differences in coefficient scale between **c** and **f** are due to minor differences in the regression procedure, specifically the use of binary indicators for reward history vs. continuous values (ranging from 0.35 to 0.6). For comparison, Tsutsui et al.⁷ reported regression results for the effect of past reward that indicated a similar temporal profile, but peaked at lag 1 with a coefficient of 1.25. Panels **a–c** adapted with permission from Tsutsui et al.⁷ and Lau and Glimcher³³.

Results

Despite its simplicity, the meta-RL framework can account for a surprising range of neuroscientific findings, including many of the results that have presented difficulties for the standard RPE model of DA. To survey the framework's ramifications, we conducted a set of six simulation experiments, each focusing on a set of experimental findings chosen to illustrate a core aspect of the theory. Throughout these simulations, we continue to apply the simple computational architecture from Fig. 1a. This approach obviously abstracts over many neuroanatomical and physiological details (see Supplementary Fig. 3). However, it allows us to demonstrate key effects in a setting where their computational origins can be readily identified. In the same spirit, our simulations focus not on parameter-dependent fits to data, but instead only on robust qualitative effects.

Simulation 1: RL in the prefrontal network.

We begin with the finding that PFC encodes recent actions and rewards, integrating these into an evolving representation of choice value. To demonstrate how meta-RL explains this finding, we simulate results from Tsutsui et al.⁷ and Lau and Glimcher³³. Both studies employed a task in which monkeys chose between two visual saccade targets, each of which yielded juice reward with a probability that intermittently changed. Given the reward schedule, an optimal strategy could be approximated by sampling each target in proportion to its yield³³, and the monkeys displayed such 'probability matching' (Fig. 2a). Recordings from dorsolateral PFC⁷ revealed activity encoding the target selected on the previous trial, the reward received, the updated values of the targets, and each upcoming action (Fig. 2b), supporting the conclusion that PFC neurons reflect the trial-by-trial construction and updating of choice value from recent experience¹⁰.

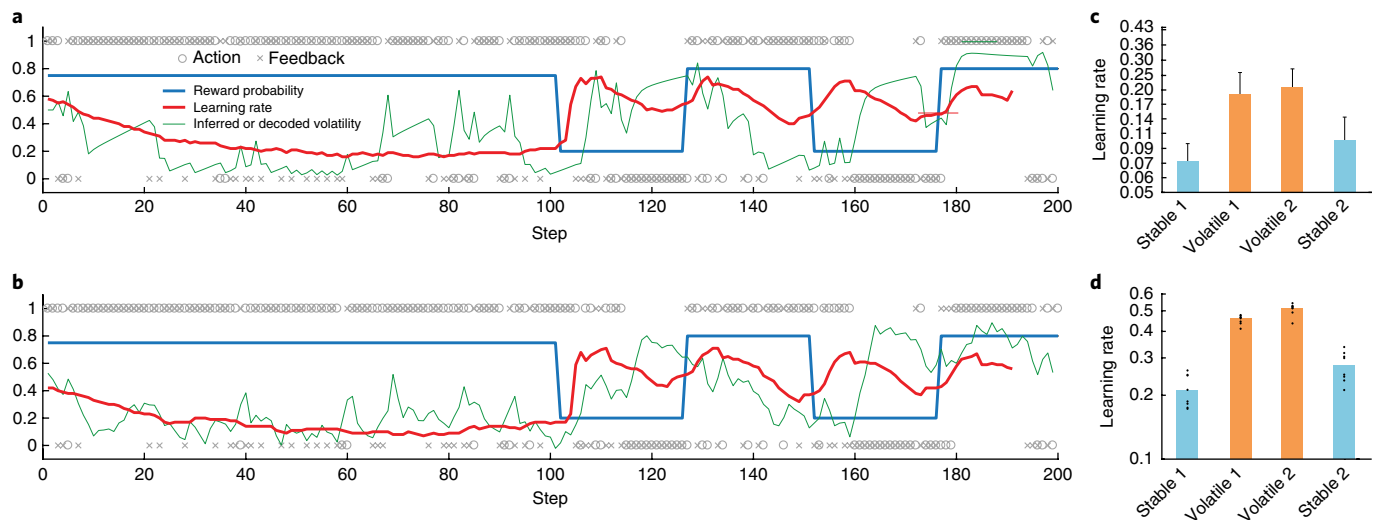


Fig. 3 | Learned RL algorithm dynamically adapts its learning rate to the volatility of the environment. **a**, Sample behavior of the Bayesian model proposed by Behrens et al.³⁴ on their volatile bandit task. Blue, true reward probability for action 1 (probabilities for actions 1 and 2 summed to 1); o, actions (action 1 above, 2 below); x, outcomes (reward above, non-reward below). Green, estimated volatility; red, learning rate. **b**, Corresponding quantities for the meta-RL model, with volatility decoded as described in Methods. **c**, Summary of human behavior redrawn from Behrens et al.³⁴, showing mean learning rates across test blocks with s.e.m. **d**, Corresponding meta-RL behavior showing a similar pattern, with higher learning rates in volatile task blocks. Dots represent data from separate training runs ($n=8$) employing different random seeds. Panels **a** and **c** adapted with permission from Behrens et al.³⁴.

We simulated these results by training our meta-RL model on the same task (see Methods). The network adapted its behavior to each task instantiation, displaying probability matching that closely paralleled experimental results (Fig. 2d). Finer-grained analyses yielded patterns similar to those observed in monkey behavior (Fig. 2c,f and Supplementary Fig. 4). Note that the network's weights were fixed at test, meaning that the network's behavior can only result from the dynamics of prefrontal recurrent neural network activity. Consistent with this, we found units in the prefrontal network that encoded each of the variables reported by Tsutsui et al.⁷ (Fig. 2e).

These results provide another basic illustration of meta-RL in action, with DA-driven training giving rise to an independent prefrontal-based learning algorithm. In this case, meta-RL accounts for both the behavior and PFC activity observed in the experimental data, but also offers an explanation for how these both may have emerged through DA-driven learning.

Simulation 2: adaptation of prefrontal-based learning to the task environment. A key aspect of meta-RL is that the learning algorithm that arises in the prefrontal network can differ from the DA-based algorithm that created it. We illustrate this principle here, focusing on differences in a single parameter: the learning rate. To this end, we simulated results from an experiment by Behrens et al.³⁴. This involved a two-armed bandit task that alternated between 'stable' periods where payoff probabilities held steady and 'volatile' periods where they fluctuated. Behrens et al.³⁴ found that human participants adopted a faster learning rate during volatile periods than stable periods, in accord with the optimal strategy identified by a Bayesian model that adjusts its learning rate on the basis of dynamic estimates of task volatility (compare Fig. 3a,c). Neuroimaging data identified a region within PFC (specifically, dorsal cingulate cortex) whose activity tracked the Bayesian model's volatility estimates.

Figure 3b,d summarizes the behavior of our model, with fixed weights, on the same task. The model was previously trained on a series of episodes with shifting volatilities (see Methods and Supplementary Fig. 5) to mimic the prior experience of different volatilities that human participants brought to the task. Like human learners, the network dynamically adapted its learning rate to the changing volatility. Moreover, as in the PFC activity observed with

fMRI, many units in the long short-term memory (LSTM) network ($37 \pm 1\%$) explicitly tracked the changing volatility.

Critically, the learning rates indicated in Fig. 3d are orders of magnitude larger than the one governing the DA-based RL algorithm that adjusted the prefrontal network's weights during training (which was set at 0.0005). The results thus provide a concrete illustration of the point that the learning algorithm that arises from meta-RL can differ from the algorithm that originally engendered it. The results also allow us to emphasize an important corollary of this principle, which is that meta-RL produces a prefrontal learning algorithm that is in fact adapted to the task environment. In the present case, this adaptation manifests in the way the learning rate responds to task fluctuations. Previous studies have proposed special-purpose mechanisms to explain dynamic shifts in learning rate^{29,35}. Meta-RL explains these shifts as an emergent effect, arising from a very general set of conditions. Moreover, dynamic learning rates constitute just one possible form of specialization. When meta-RL occurs in environments with different structures, qualitatively different learning rules will emerge, a point that will be illustrated by subsequent simulations.

Simulation 3: reward prediction errors reflecting inferred value.

As introduced earlier, one apparent challenge for the standard model is that the DA RPE signal reflects knowledge of task structure. An example is the 'inferred-value' effect reported by Bromberg-Martin et al.¹². In each trial in their task, a visual target appeared on either the left or right side of a display, and the monkey was expected to saccade to that target. At any point in the experiment, either the left or right target yielded a juice reward while the other target did not, and these role assignments reversed intermittently throughout the testing session. The key observation concerned DA signals following a reversal: immediately after the monkey experienced a change in reward for one target, the DA response to the appearance of the other target changed dramatically, reflecting an inference that the value of that target had also changed (Fig. 4a).

This inferred-value effect, along with other related findings, has given rise to models in which either PFC or hippocampus encodes abstract latent-state representations^{19,36–38}, which can then feed into the computations generating the RPE¹³. As it turns out, a related

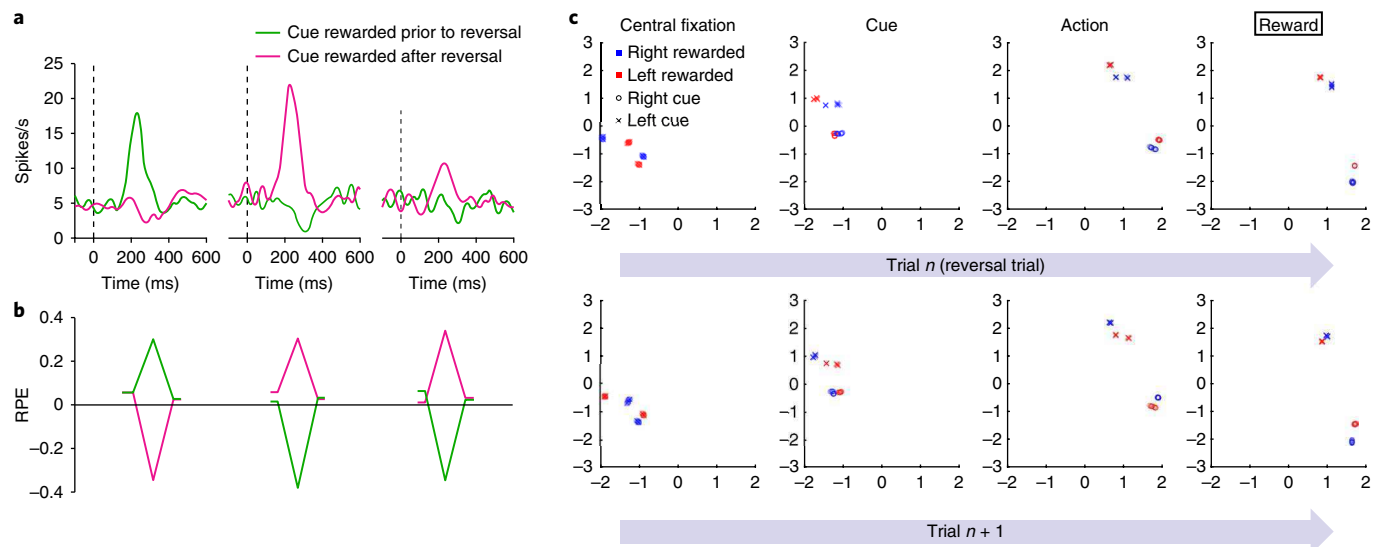


Fig. 4 | Reward prediction errors of the simulated agent reflect inferred value, not just experienced value, similarly to what is observed in monkeys.

a. Results adapted with permission from Bromberg-Martin et al.¹². Dopaminergic activity in response to cues ahead of a reversal (left) and for cues with an experienced (middle) and inferred (right) change in value. Curves represent average neuronal responses ($n=42\text{--}63$ for each curve) from 2 monkeys. **b.** Corresponding RPE signals from the model. Leading and trailing points for each data series correspond to initial fixation and saccade steps. Peaks and troughs correspond to stimulus presentation. Although the RPE responses in empirical data are smaller for inferred value, which is not shown here, other empirical results from Bromberg-Martin et al.¹² depict much more similar patterns. In particular, recordings from lateral habenula in the same task, but from different animals, yielded much more similar RPE signals (and behavioral reaction times) across conditions. **c.** First two principal components of recurrent neural network activity (LSTM output) on individual steps in the task from simulation 3, analyzed from 1,200 evaluation episodes performed by 1 trained network replica. Other replicas yielded very similar results. The top row focuses on reversal trials, where the final reward feedback (rightmost panel) signals that the rewarded target has switched relative to preceding trials. The second row focuses on trials immediately following the reversal trials examined in the first row. Activity patterns cluster according to the current latent state of the task (i.e., which cue is currently rewarded), and later in each trial also according to the action selected. As shown in the leftmost panels ('central fixation'), at the beginning of each trial, the network's activation state represents the latent state of the task, and this abruptly reverses following reversal trials.

mechanism arises naturally from meta-RL. To show this, we trained our meta-RL model on the task of Bromberg-Martin et al.¹². At test, we observed RPE signals that reproduced the pattern displayed by DA. In particular, the model clearly reproduced the critical inferred-value effect (Fig. 4b).

The explanation for this result is straightforward. In our architecture, the reward prediction component of the RPE comes from the prefrontal network's state-value output, consistent with data showing that DA signaling is influenced by projections from PFC¹³. As DA-based training causes the prefrontal network to encode information about task dynamics (Fig. 4c), this information also manifests in the reward prediction component of the RPE.

Simulation 4: 'model-based' behavior—the two-step task.

Another important setting where structure-sensitive DA signaling has been observed is in tasks designed to probe for model-based control. Perhaps the most heavily studied task in this category is the 'two-step' model introduced by Daw et al.¹⁵. In the version we will consider (Fig. 5a), every trial begins with a decision between two actions. Each triggers a probabilistic transition to one of two perceptually distinguishable second-stage states; for each action there is a 'common' (high probability) transition and an 'uncommon' (low probability) one. Each second-stage state then delivers a reward with a particular probability, which changes intermittently (see Methods).

The interest of the two-step task lies in its ability to differentiate between model-free and model-based learning. Most important is behavior following an uncommon transition. Consider an action selected at stage 1 that triggers an uncommon transition followed by a reward. Model-free learning, driven by action-reward associations, will increase the probability of repeating the same first-stage

action on the subsequent trial. Model-based learning, in contrast, will take account of the task's transition structure and increase the probability of choosing the opposite action (Fig. 5b). Both humans¹⁵ and rodents³⁹ typically display behavioral evidence of model-based learning (Fig. 5c,d).

We trained our meta-RL model on the two-step task (see Methods) and found that its behavior at test assumed the form associated with model-based control (Fig. 5c,e). As in previous simulations, the network was tested with weights frozen. The observed behavior was thus generated by the dynamics of the recurrent prefrontal network. However, it is important to recall that the network was trained by a model-free DA-driven RL algorithm. In a neuroscientific context, this points to the interesting possibility that the PFC's implementation of model-based RL might in fact arise from model-free DA-driven training (see Supplementary Figs. 6–8 for further analysis).

Critically, the two-step task was one of the first contexts in which structure-sensitive RPE signals were reported. In particular, using fMRI, Daw et al.¹⁵ observed RPEs in human ventral striatum—a major target for DA—that tracked those predicted by a model-based RL algorithm. The same effect arises in our meta-RL model (Fig. 5f,g).

Simulation 5: learning to learn. Simulations 3 and 4 focused on scenarios involving an alternation between two versions of a task, each of which might become familiar over time. Here we apply meta-RL to a task in which new stimuli are continually presented, requiring learning in the fullest sense. In this setting, we show that meta-RL can account for situations in which past experience speeds new learning, an effect often called 'learning to learn'.

In the task that originally inspired this term, Harlow⁴⁰ presented a monkey with two unfamiliar objects, one covering a well contain-

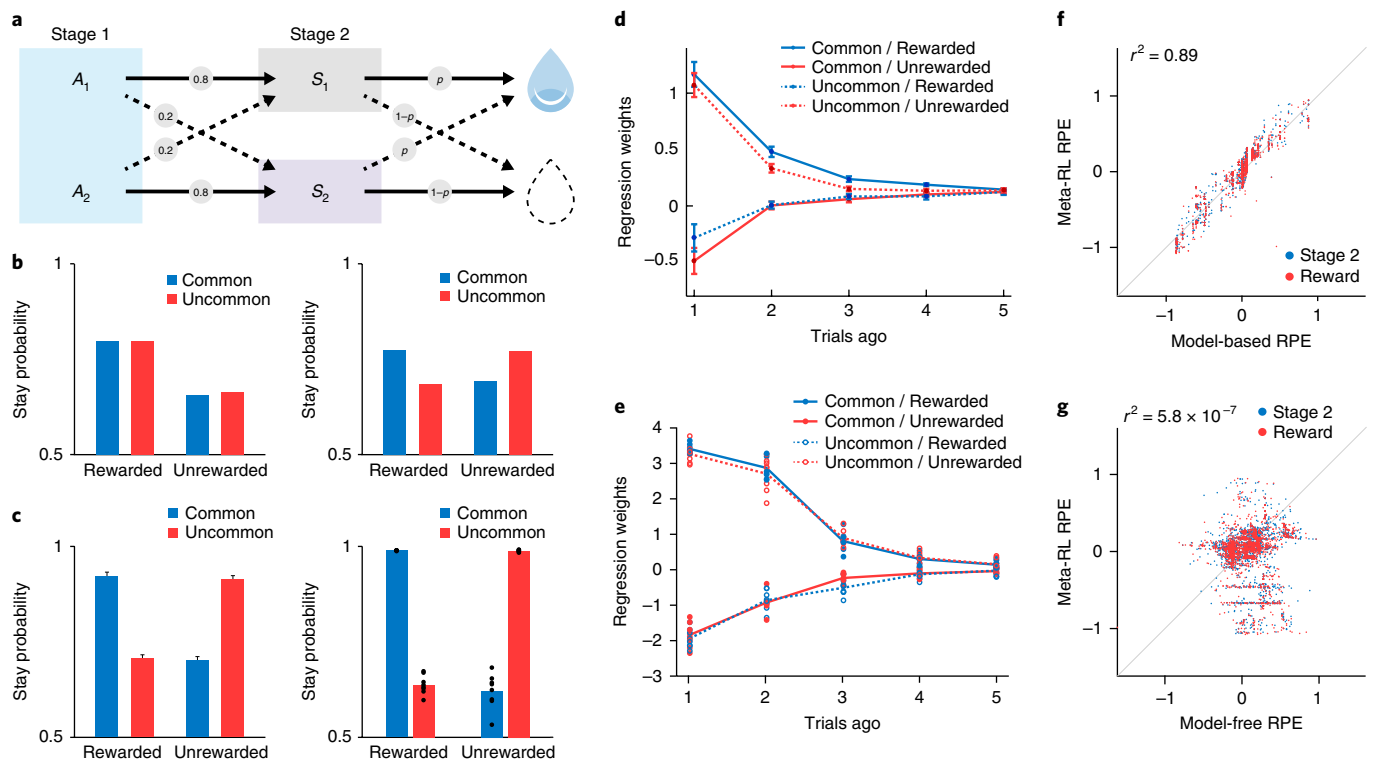


Fig. 5 | Learned RL algorithm displays model-based behavior and RPEs. **a**, Structure of the two-step task^{15,39}, depicting how first-stage actions A_1 and A_2 probabilistically transition to second-stage states S_1 and S_2 . **b**, Top, canonical pattern of behavior for model-free (left) and model-based (right) learning. **c**, Left, data from Miller et al.³⁹ showing a canonical model-based pattern in animal behavior. Plot represents mean stay probability for $n = 6$ rats; error bars are s.e.m. For full distribution see Miller et al.³⁹. Right, corresponding behavior of the network. Dots represent data from separate training runs using different random seeds ($n = 8$). **d**, More detailed analysis from Miller et al.³⁹, showing the influence of transition types and rewards, at multiple trial lags, on animal behavior. Plot represents mean regression weights for $n = 21$ rats; error bars are s.e.m. The pattern rules out an alternative, model-free mechanism (see Supplementary Fig. 6). **e**, Corresponding analysis of meta-RL behavior. Dots represent data from separate training runs using different random seeds ($n = 8$). **f, g**, Regression of model RPE signals against predictions from model-based (**f**) and model-free (**g**) algorithms at different stages of the trial. Analysis done for 30 evaluation episodes, each consisting of 100 trials, performed by 1 fully trained network out of 8 replicas. Panels **c** (left) and **d** adapted with permission from Miller et al.³⁹.

ing food reward, the other an empty well. The animal chose freely between the objects and could retrieve the food reward if present. The left–right positions of the objects were then randomly reset, and a new trial began. After six repetitions of this process, two entirely new objects were substituted, and the process began again. Within each block of trials, one object was chosen to be consistently rewarded, with the other consistently unrewarded. Early in training, monkeys were slow to converge on the correct object in each block. But after substantial practice, monkeys showed perfect performance after only a single trial, reflecting an understanding of the task's rules (Fig. 6b).

To evaluate the ability of meta-RL to account for Harlow's⁴⁰ results, we converted his task to one requiring choices between images presented on a simulated computer display (Fig. 6a and Supplementary Video 1). The task was otherwise unchanged, with a pair of unfamiliar images introduced after every six trials. To allow our model to process high-dimensional pixel inputs, we augmented it with a conventional image-processing network (see Methods). The resulting system generated learning curves closely resembling the empirical data (Fig. 6c). Following training, the network learned in a single trial how to respond in each new block, replicating Harlow's⁴⁰ learning-to-learn effect.

Simulation 6: the role of dopamine—effects of optogenetic manipulation. The meta-RL framework requires that the inputs to the prefrontal network contain information about recent rewards.

Our implementation satisfies this requirement by feeding in a scalar signal explicitly representing the amount of reward received on the previous time-step (Fig. 1a). However, any signal that is robustly correlated with reward would suffice. Sensory signals, such as the taste of juice, could thus play the requisite role. Another particularly interesting possibility is that information about rewards might be conveyed by DA itself. In this case, DA would play two distinct roles. First, as in the standard model, DA would modulate synaptic plasticity in the prefrontal network. Second, the same DA signal would support activity-based RL computations in the prefrontal network by injecting information about recent rewards⁴¹. We modified our original network to provide the RPE, in place of reward, as input to the network (see Methods and Supplementary Fig. 9) and found that it generated comparable behavior on the tasks from simulations 1–5.

This variation on meta-RL suggests a new interpretation for the findings of recent experiments using optogenetic techniques to block or induce dopaminergic RPE signals^{42,43}. To illustrate, we consider an experiment by Stopper et al.⁴⁴ involving a two-armed bandit task with intermittently shifting payoff probabilities (see Methods). Blocking DA activity during delivery of food rewards from one lever led to a reduced preference for that lever. Conversely, artificially stimulating DA release when one lever failed to yield food increased preference for that lever (Fig. 7a). We simulated these conditions by incrementing or decrementing the RPE input to the network in Supplementary Fig. 9 and observed comparable behavioral results

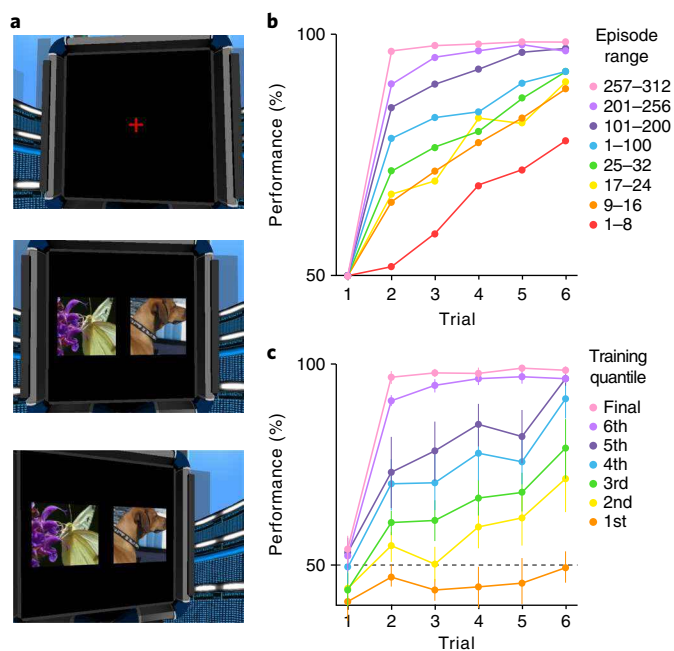


Fig. 6 | Meta-RL learns to learn about abstract structure and novel stimuli in a visually rich 3D environment. **a**, Example image inputs for simulation 5 (full resolution), showing fixation cross (top), initial stimulus presentation (middle) and saccade outcome (bottom). **b**, Accuracy data from Harlow⁴⁰ for each step following introduction of new objects and at multiple points in training of connection weights. **c**, Model performance at seven successive stages of training (see Methods), averaged across all networks ($n=14$) that achieved maximal performance (out of 50 replicas). Error bars represent 95% confidence intervals.

(Fig. 7b). As in previous simulations, these results were obtained with the network weights held constant. The shift in behavior induced by our simulated optogenetic intervention thus reflects the impact of the dopaminergic RPE signal on unit activities within the prefrontal network, rather than an effect on synaptic weights. In this regard, the present simulation thus provides an interpretation of the experimental results that is radically different from the one that would proceed from the standard model.

Functional neuroanatomy. As our simulation results show, meta-RL provides a new lens through which to examine the respective functions of the prefrontal network and the DA system, and their interactions during learning. As we have noted, a key assumption of the meta-RL theory is that the prefrontal network can be understood as a recurrent neural circuit^{16–18}. In fact, the PFC lies at the intersection of multiple recurrent circuits, involving recurrent connections within PFC itself, connections with other cortical regions⁴⁵ and, most crucially, connections with dorsal striatum and medio-dorsal thalamus (Supplementary Fig. 3). A prevalent view¹⁷ is that the role of the striatum within this cortico–basal ganglia–thalamo–cortical loop is to regulate the dynamics of PFC by gating the flow of information into PFC circuits. According to this account, DA serves as a training signal within the striatum, shaping the operation of the gating function through temporal difference learning. Replicating our simulation results using this more detailed model of the cortico–striatal loop represents a worthwhile target for future research.

It is worth noting that the striatal gating theory was originally inspired by LSTM networks, of the kind we have employed in our simulations¹⁷ (Fig. 1b). Indeed, recent work posits multiple gating mechanisms operating on PFC, which directly parallel the input, memory and output gating mechanisms traditionally implemented

in LSTM networks⁴⁶ (see Methods). While these parallels are intriguing, it should also be noted that other recent work has modeled RL in the PFC using more generic recurrent neural networks that lack the gating mechanisms involved in LSTM networks^{16,18}. Although the meta-RL effect we have identified relies only on the presence of recurrence, and not on any particular gating mechanism, it would be interesting to test whether particular gating mechanisms provide a more precise quantitative fit to human and animal learning curves and neural representations.

The neuroscientific data we have considered in this paper center mostly on structures in the so-called ‘associative loop’ connecting the dorsal PFC with the basal ganglia (Supplementary Fig. 3). However, this is only one of several such recurrent loops in the brain, with others running through other sectors of cortex (Supplementary Fig. 3). This raises the question of whether other recurrent loops, such as the ‘sensorimotor loop’ running through somatosensory and motor cortices, also support meta-RL. Although we cannot rule out the possibility that forms of meta-learning may also emerge in these loops, we observe that the meta-RL effect described in this paper emerges only when network inputs carry information about recent actions and rewards and when network dynamics support maintenance of information over suitable time periods. These two factors may differentiate the associative loop from other cortico–basal ganglia–thalamo–cortical loops⁴⁷ explaining why meta-RL might arise uniquely from this circuit (Supplementary Fig. 3).

Most of our simulations modeled all of PFC as a single fully connected network without regional specialization. However, there are important functional-anatomical distinctions within PFC. For example, the probability-matching behavior seen in simulation 1 and the model-based learning pattern in simulation 4 have robust neural correlates in dorsolateral PFC^{7,11}, while the volatility coding modeled in simulation 2 was reported in the anterior cingulate cortex³⁴. Although the differential roles of anterior cingulate and dorsolateral PFC are still under active debate, an important next step for the present theory will be to consider how the computations involved in meta-RL might play out across these regions, given their different cell properties, internal circuitry and extrinsic connectivity. Also relevant are PFC regions lying in the so-called ‘limbic loop’, including orbitofrontal and ventromedial PFC (see Supplementary Fig. 3). Both of the latter regions have been implicated in reward coding, while recent work suggests that the orbitofrontal cortex may additionally encode abstract latent states^{36,37}—both critical functions for meta-RL, as illustrated in several of our simulations. Once again, a fuller development of the meta-RL theory will need to incorporate the relative roles of these regions more explicitly.

Discussion

We have put forward a new proposal concerning the roles of DA and PFC in reward-based learning, leveraging the notion of meta-RL. The framework we have advanced conserves the standard RPE model of DA function but places it in a new context, allowing it to newly accommodate previously puzzling findings. As our simulations have shown, meta-RL accounts for a diverse range of observations concerning both DA and PFC function, providing a bridge between the literatures addressing these two systems.

In addition to explaining existing data, the meta-RL framework also leads to a number of testable predictions. As we have seen, the theory suggests that the role of PFC in model-based control could arise, at least in part, from DA-driven synaptic learning. If this is correct, then interfering with phasic DA signaling during initial training should interfere with the emergence of model-based control in tasks like those studied in simulations 4 and 5. Another prediction concerns model-based dopaminergic RPE signals. Meta-RL attributes these signals to value inputs from the prefrontal network. If this is correct, then lesioning or inactivating PFC or its associated striatal nuclei should eliminate model-based

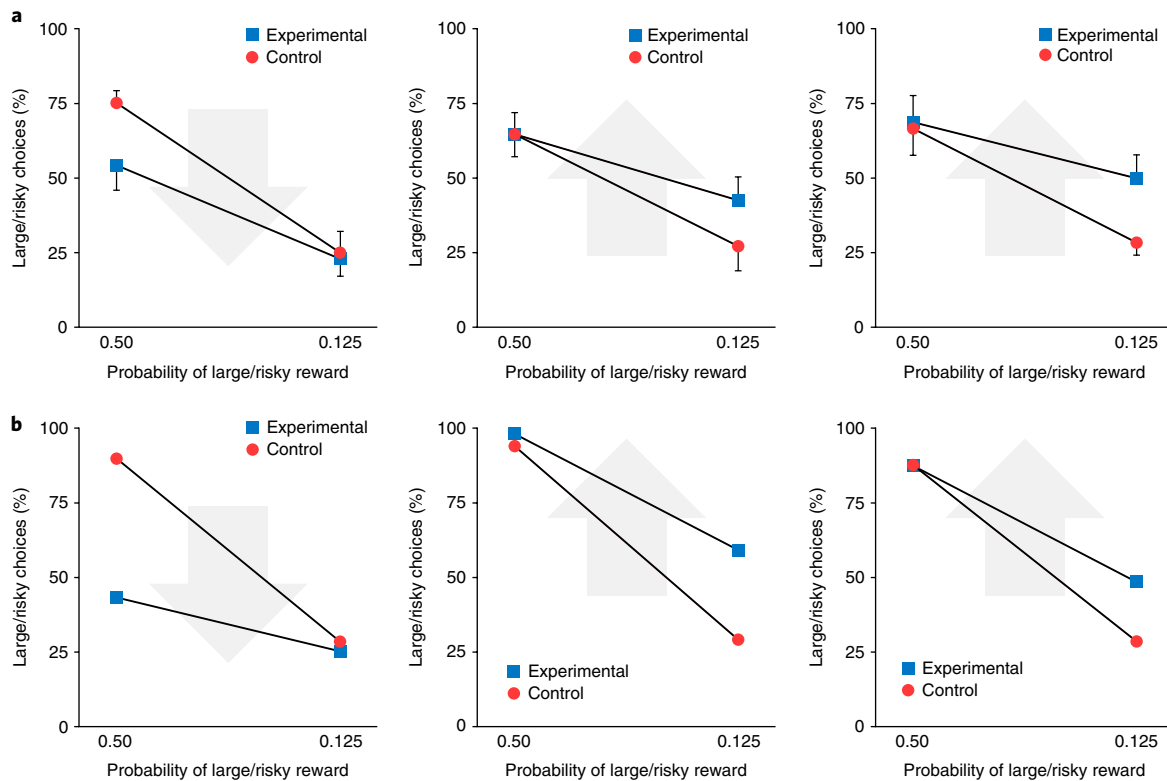


Fig. 7 | Dual roles of dopamine: producing synaptic change for slow learning and carrying information about rewards for fast learning. **a**, Behavioral data adapted with permission from Stopper et al.⁴⁴. Their task posed a recurring choice between a lever that reliably yielded a small reward (small/certain) and one that yielded a large reward with probability p (large/risky). In a control condition, rats chose the large/risky lever more often when $p=0.5$ than when $p=0.125$ (blue). Left, effect of optogenetically blocking DA when large/risky rewards occurred ($n=11$). Arrow indicates direction of shift with optogenetic intervention. Center, effect of blocking DA when small/certain rewards occurred ($n=8$). Right, effect of triggering DA when large/risky rewards failed to occur ($n=9$). All plots represent mean \pm s.e.m. **b**, Model behavior in circumstances modeling those imposed by Stopper et al.⁴⁴, with conditions presented as in **a**.

DA signaling⁴⁸. Further predictions might be specified by examining the patterns of activity arising in the prefrontal network portion of our model, treating these as predictors for neural activity in animals performing relevant tasks. Behavioral tasks useful to such an undertaking might be drawn from recent work implicating PFC in the identification of latent states^{19,36–38} and abstract rules⁴⁹ in RL contexts.

Meta-RL also raises a range of broader questions that we hope will stimulate new experimental work. What might be the relative roles of mesolimbic, mesocortical and nigrostriatal DA pathways in a meta-RL context? Does meta-RL point to new interpretations of the division of labor between dorsal and ventral, or medial and lateral, sectors of PFC? How should we interpret data pointing to the existence of systems supporting both model-based and model-free RL (Supplementary Figs. 7 and 8)? What new dynamics emerge when meta-RL is placed in contact with mechanisms subserving episodic memory⁵⁰? All in all, meta-RL offers a new orienting point in the landscape of ideas concerning reward-based learning, one that may prove useful in the development of new research questions and in the interpretation of new findings.

Methods

Methods, including statements of data availability and any associated accession codes and references, are available at <https://doi.org/10.1038/s41593-018-0147-8>.

Received: 13 December 2017; Accepted: 5 April 2018;

Published online: 14 May 2018

References

- Sutton, R. S. & Barto, A. G. *Reinforcement Learning: An Introduction* (MIT Press, Cambridge, MA, USA, 1998).
- Montague, P. R., Dayan, P. & Sejnowski, T. J. A framework for mesencephalic dopamine systems based on predictive Hebbian learning. *J. Neurosci.* **16**, 1936–1947 (1996).
- Daw, N. D. & Tobler, P. N. Value learning through reinforcement: the basics of dopamine and reinforcement learning. *Neuroeconomics: Decision Making and the Brain* 2nd edn. (eds. Glimcher, P. W. & Fehr, E.) 283–298 (Academic, New York, 2014).
- Rushworth, M. F. & Behrens, T. E. Choice, uncertainty and value in prefrontal and cingulate cortex. *Nat. Neurosci.* **11**, 389–397 (2008).
- Seo, H. & Lee, D. Cortical mechanisms for reinforcement learning in competitive games. *Phil. Trans. R. Soc. Lond. B* **363**, 3845–3857 (2008).
- Padoa-Schioppa, C. & Assad, J. A. Neurons in the orbitofrontal cortex encode economic value. *Nature* **441**, 223–226 (2006).
- Tsutsui, K., Grabenhorst, F., Kobayashi, S. & Schultz, W. A dynamic code for economic object valuation in prefrontal cortex neurons. *Nat. Commun.* **7**, 12554 (2016).
- Kim, J.-N. & Shadlen, M. N. Neural correlates of a decision in the dorsolateral prefrontal cortex of the macaque. *Nat. Neurosci.* **2**, 176–185 (1999).
- Seo, M., Lee, E. & Averbeck, B. B. Action selection and action value in frontal-striatal circuits. *Neuron* **74**, 947–960 (2012).
- Barraclough, D. J., Conroy, M. L. & Lee, D. Prefrontal cortex and decision making in a mixed-strategy game. *Nat. Neurosci.* **7**, 404–410 (2004).
- Daw, N. D., Niv, Y. & Dayan, P. Uncertainty-based competition between prefrontal and dorsolateral striatal systems for behavioral control. *Nat. Neurosci.* **8**, 1704–1711 (2005).
- Bromberg-Martin, E. S., Matsumoto, M., Hong, S. & Hikosaka, O. A pallidum-habenula-dopamine pathway signals inferred stimulus values. *J. Neurophysiol.* **104**, 1068–1076 (2010).
- Nakahara, H. & Hikosaka, O. Learning to represent reward structure: a key to adapting to complex environments. *Neurosci. Res.* **74**, 177–183 (2012).

14. Sadacca, B. F., Jones, J. L. & Schoenbaum, G. Midbrain dopamine neurons compute inferred and cached value prediction errors in a common framework. *Elife* **5**, e13665 (2016).
15. Daw, N. D., Gershman, S. J., Seymour, B., Dayan, P. & Dolan, R. J. Model-based influences on humans' choices and striatal prediction errors. *Neuron* **69**, 1204–1215 (2011).
16. Mante, V., Sussillo, D., Shenoy, K. V. & Newsome, W. T. Context-dependent computation by recurrent dynamics in prefrontal cortex. *Nature* **503**, 78–84 (2013).
17. O'Reilly, R. C. & Frank, M. J. Making working memory work: a computational model of learning in the prefrontal cortex and basal ganglia. *Neural Comput.* **18**, 283–328 (2006).
18. Song, H. F., Yang, G. R. & Wang, X.-J. Reward-based training of recurrent neural networks for cognitive and value-based tasks. *Elife* **6**, e21492 (2017).
19. Redish, A. D., Jensen, S., Johnson, A. & Kurth-Nelson, Z. Reconciling reinforcement learning models with behavioral extinction and renewal: implications for addiction, relapse, and problem gambling. *Psychol. Rev.* **114**, 784–805 (2007).
20. Haber, S. N. The place of dopamine in the cortico-basal ganglia circuit. *Neuroscience* **282**, 248–257 (2014).
21. Frank, M. J., Seeberger, L. C. & O'Reilly, R. C. By carrot or by stick: cognitive reinforcement learning in parkinsonism. *Science* **306**, 1940–1943 (2004).
22. Houk, J. C., Adams, C. M. & Barto, A. G. A model of how the basal ganglia generate and use neural signals that predict reinforcement. in *Models of Information Processing in the Basal Ganglia* (eds. Houk, J.C. & Davis, D.G.) 249–270 (MIT Press, Cambridge, MA, USA, 1995).
23. Rougier, N. P., Noelle, D. C., Braver, T. S., Cohen, J. D. & O'Reilly, R. C. Prefrontal cortex and flexible cognitive control: rules without symbols. *Proc. Natl. Acad. Sci. USA* **102**, 7338–7343 (2005).
24. Acuna, D. E. & Schrater, P. Structure learning in human sequential decision-making. *PLoS Comput. Biol.* **6**, e1001003 (2010).
25. Collins, A. G. & Frank, M. J. How much of reinforcement learning is working memory, not reinforcement learning? A behavioral, computational, and neurogenetic analysis. *Eur. J. Neurosci.* **35**, 1024–1035 (2012).
26. Thrun, S. & Pratt, L. *Learning to Learn* (Springer Science & Business Media, New York, 2012).
27. Khamassi, M., Enel, P., Dominey, P. F. & Procyk, E. Medial prefrontal cortex and the adaptive regulation of reinforcement learning parameters. *Prog. Brain Res.* **202**, 441–464 (2013).
28. Ishii, S., Yoshida, W. & Yoshimoto, J. Control of exploitation-exploration meta-parameter in reinforcement learning. *Neural Netw.* **15**, 665–687 (2002).
29. Schweighofer, N. & Doya, K. Meta-learning in reinforcement learning. *Neural Netw.* **16**, 5–9 (2003).
30. Schmidhuber, J., Zhao, J. & Wiering, M. Simple principles of metalearning. IDISA (Istituto Dalle Molle Di Studi Sull Intelligenza Artificiale) Technical Report **69-96**, 1–23 (1996).
31. Wang, J.X. et al. Learning to reinforcement learn. Preprint at <https://arxiv.org/abs/1611.05763> (2016).
32. Duan, Y. et al. RL2: fast reinforcement learning via slow reinforcement learning. Preprint at <https://arxiv.org/abs/1611.02779> (2016).
33. Lau, B. & Glimcher, P. W. Dynamic response-by-response models of matching behavior in rhesus monkeys. *J. Exp. Anal. Behav.* **84**, 555–579 (2005).
34. Behrens, T. E. J., Woolrich, M. W., Walton, M. E. & Rushworth, M. F. S. Learning the value of information in an uncertain world. *Nat. Neurosci.* **10**, 1214–1221 (2007).
35. Iigaya, K. Adaptive learning and decision-making under uncertainty by metaplastic synapses guided by a surprise detection system. *Elife* **5**, e18073 (2016).
36. Schuck, N. W., Cai, M. B., Wilson, R. C. & Niv, Y. Human orbitofrontal cortex represents a cognitive map of state space. *Neuron* **91**, 1402–1412 (2016).
37. Chan, S. C., Niv, Y. & Norman, K. A. A probability distribution over latent causes, in the orbitofrontal cortex. *J. Neurosci.* **36**, 7817–7828 (2016).
38. Hampton, A. N., Bossaerts, P. & O'Doherty, J. P. The role of the ventromedial prefrontal cortex in abstract state-based inference during decision making in humans. *J. Neurosci.* **26**, 8360–8367 (2006).
39. Miller, K. J., Botvinick, M. M. & Brody, C. D. Dorsal hippocampus contributes to model-based planning. *Nat. Neurosci.* **20**, 1269–1276 (2017).
40. Harlow, H. F. The formation of learning sets. *Psychol. Rev.* **56**, 51–65 (1949).
41. Trujillo-Pisanty, I., Solis, P., Conover, K., Dayan, P. & Shizgal, P. On the forms of learning supported by rewarding optical stimulation of dopamine neurons. *Soc. Neurosci. Annu. Meet.* 66.06, <http://www.abstractsonline.com/pp8/#!/4071/presentation/29633> (2016).
42. Kim, K. M. et al. Optogenetic mimicry of the transient activation of dopamine neurons by natural reward is sufficient for operant reinforcement. *PLoS One* **7**, e33612 (2012).
43. Chang, C. Y. et al. Brief optogenetic inhibition of dopamine neurons mimics endogenous negative reward prediction errors. *Nat. Neurosci.* **19**, 111–116 (2016).
44. Stopper, C. M., Tse, M. T. L., Montes, D. R., Wiedman, C. R. & Floresco, S. B. Overriding phasic dopamine signals redirects action selection during risk/reward decision making. *Neuron* **84**, 177–189 (2014).
45. Wang, X.-J. Synaptic reverberation underlying mnemonic persistent activity. *Trends Neurosci.* **24**, 455–463 (2001).
46. Chatham, C. H. & Badre, D. Multiple gates on working memory. *Curr. Opin. Behav. Sci.* **1**, 23–31 (2015).
47. Kim, H., Lee, D. & Jung, M. W. Signals for previous goal choice persist in the dorsomedial, but not dorsolateral striatum of rats. *J. Neurosci.* **33**, 52–63 (2013).
48. Takahashi, Y. K. et al. Expectancy-related changes in firing of dopamine neurons depend on orbitofrontal cortex. *Nat. Neurosci.* **14**, 1590–1597 (2011).
49. Collins, A. G. E. & Frank, M. J. Neural signature of hierarchically structured expectations predicts clustering and transfer of rule sets in reinforcement learning. *Cognition* **152**, 160–169 (2016).
50. Gershman, S. J. & Daw, N. D. Reinforcement learning and episodic memory in humans and animals: An integrative framework. *Annu. Rev. Psychol.* **68**, 101–128 (2017).

Acknowledgements

We are grateful to K. Miller, F. Grabenhorst, T. Behrens, E. Bromberg-Martin, S. Floresco and P. Glimcher for graciously providing help with and permission for adapting their data. We thank C. Blundell and R. Munos for discussions and comments on an earlier draft.

Author contributions

J.X.W., Z.K.-N., and M.B. designed the simulations. J.X.W. and Z.K.-N. performed the simulations and analyzed the data. D.T., H.S. and J.Z.L. contributed and helped with code. All authors wrote the manuscript.

Competing interests

The authors are employed by DeepMind Technologies Limited.

Additional information

Supplementary information is available for this paper at <https://doi.org/10.1038/s41593-018-0147-8>.

Reprints and permissions information is available at www.nature.com/reprints.

Correspondence and requests for materials should be addressed to M.B.

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

Methods

Architecture and learning algorithm. All of our simulations employed a common set of methods, with minor implementational variations. The agent architecture centers on a fully connected, gated recurrent neural network (LSTM: long short-term memory network⁵¹; see below for equations). In all experiments except where specified, the input included the observation, a scalar indicating the reward received on the preceding time-step, and a one-hot representation of the action taken on the preceding time-step. The outputs consisted of a scalar baseline (value function) and a real vector with length equal to the number of available actions. Actions were sampled from the softmax distribution defined by this vector. Some other architectural details were varied as required by the structure of different tasks (see simulation-specific details below). Reinforcement learning was implemented by Advantage Actor-Critic, as detailed by Mnih et al.⁵². Details of training, including the use of entropy regularization and a combined policy and value estimate loss, are described by Mnih et al.⁵². In brief, the gradient of the full objective function is the weighted sum of the policy gradient, the gradient with respect to the state-value function loss, and an entropy regularization term, defined as follows:

$$\begin{aligned}\nabla \mathcal{L} &= \nabla \mathcal{L}_\pi + \nabla \mathcal{L}_v + \nabla \mathcal{L}_{\text{ent}} \\ &= \frac{\partial \log \pi(a_t | s_t; \theta)}{\partial \theta} \delta_t(s_t; \theta_v) + \beta_v \delta_t(s_t; \theta_v) \frac{\partial V}{\partial \theta_v} \\ &\quad + \beta_e \left[\frac{\partial H(\pi(a_t | s_t; \theta))}{\partial \theta} \right] \\ \delta_t(s_t; \theta_v) &= [R_t - V(s_t; \theta_v)] \\ R_t &= \sum_{i=0}^{k-1} [\gamma^i r_{t+i} + \gamma^k V(s_{t+k}; \theta_v)]\end{aligned}$$

where a_t , s_t , and R_t define the action, state and discounted n -step bootstrapped return at time t (with discount factor γ), k is the number of steps until the next terminal state and is upper bounded by the maximum unroll length t_{max} , π is the policy (parameterized by neural network parameters θ), V is the value function (parameterized by θ_v), estimating the expected return from state s , $H(\pi)$ is the entropy of the policy, and β_v and β_e are hyperparameters controlling the relative contributions of the value estimate loss and entropy regularization term, respectively. $\delta_t(s_t; \theta_v)$ is the n -step return temporal-difference error that provides an estimate of the advantage function for actor-critic. The parameters of the neural network were updated via gradient descent and backpropagation through time, using Advantage Actor-Critic as detailed by Mnih et al.⁵². Note that while the parameters θ and θ_v are being shown as separate, as in Mnih et al.⁵², in practice they share all non-output layers and differ only in the softmax output for the policy and one linear output for the value function. Simulations 1–4 and 6 used a single thread and received discrete observations coded as one-hot vectors, with length of the number of possible states (see Supplementary Fig. 1 for pseudocode for single-threaded Advantage Actor-Critic). Simulation 5 used 32 asynchronous threads during training and received RGB frames as input (see Mnih et al.⁵² for asynchronous multi-threaded algorithm and pseudocode). The core recurrent network consisted of 48 LSTM units in simulations 1–4, 256 units in simulation 5, and two separate LSTMs of 48 units each (for policy and value) in simulation 6 (see Supplementary Fig. 1).

In a standard nongated recurrent neural network, the state at time step t is a linear projection of the state at time step $t-1$, followed by a nonlinearity. This kind of ‘vanilla’ recurrent neural network can have difficulty with long-range temporal dependencies because it has to learn a very precise mapping just to copy information unchanged from one time state to the next. An LSTM, by contrast, works by copying its internal state (called the ‘cell state’) from each time step to the next. Rather than having to learn how to remember, it remembers by default. However, it is also able to choose to forget, using a ‘forget’ (or maintenance) gate, and to choose to allow new information to enter, using an ‘input’ gate. Because it may not want to output its entire memory contents at each time step, there is also an ‘output’ gate to control what to output. Each of these gates are modulated by a learned function of the state of the network.

More precisely, the dynamics of the LSTM were governed by standard equations^{51,53}:

$$\begin{aligned}i_t &= \sigma(W_{xi} x_t + W_{hi} h_{t-1} + b_i) \\ f_t &= \sigma(W_{xf} x_t + W_{hf} h_{t-1} + b_f) \\ c_t &= f_t \circ c_{t-1} + i_t \circ \tanh(W_{xc} x_t + W_{hc} h_{t-1} + b_c) \\ o_t &= \sigma(W_{xo} x_t + W_{ho} h_{t-1} + b_o)\end{aligned}$$

$$h_t = o_t \circ \tanh(c_t)$$

where x_t is the input to the LSTM at time t , h_t is the hidden state, i_t is the input gate, f_t is the forget/maintenance gate, o_t is the output gate, c_t is the cell state, σ is the sigmoid function, and \circ is an operator denoting element-wise multiplication.

General task structure. Tasks were episodic, comprising a set number of trials (fixed to a constant number per episode unless otherwise specified), with task parameters randomly drawn from a distribution and fixed for the duration of the episode. In simulations 1, 3, 4 and 6, each trial began with a fixation cue, requiring a distinct central fixation response (a_c), followed by one or more stimulus cues, each requiring a stimulus response (either left, a_l , or right, a_r). Failure to produce a valid response to either the fixation or stimulus cues resulted in a reward of -1 . In simulation 5, which involved high-dimensional visual inputs, fixation and image selection required emitting left–right actions continuously to shift the target to within the center of the field of view. An additional no-op action was provided to allow the agent to maintain fixation as necessary, and no negative reward was given for producing invalid actions. See simulation-specific methods for more details.

Training and testing. Both training and testing environments involved sampling a task from predetermined task distributions—in most cases sampling randomly, although see simulations 1 and 2 for principled exceptions—with the LSTM hidden state initialized at the beginning of each episode (initial state learned for simulations 1–4 and 6; initialized to 0 for simulation 5). Unless otherwise noted, training hyperparameters (as defined by Mnih et al.⁵²) were as follows: learning rate = 0.0007, discount factor = 0.9, state-value estimate cost $\beta_v = 0.05$, and entropy cost $\beta_e = 0.05$. Weights were optimized using Shared RMSProp and backpropagation through time⁵², which involved unrolling the recurrent network a fixed number of time-steps that ranged from $t_{\text{max}} = 100$ to 300, depending on task, and determined the number of steps when calculating the bootstrapped n -step return. The agent was then evaluated on a testing episodes, during which all network weights were held fixed. No parameter optimization was undertaken to improve fits to data, beyond the selection of what appeared to be sensible a priori values, based on prior experience with related work, and minor adjustment simply to obtain robust task acquisition.

We now detail the simulation-specific task designs, hyperparameters, and analyses. Note that details of the simulations reported in the introductory two-armed bandit simulations were drawn from Wang et al.³¹.

Simulation 1: Reinforcement learning in the prefrontal network. This task, from Tsutsui et al.⁷, required the agent to select between two actions, a_l and a_r , with dynamic reward probabilities that changed in response to previous choices. Reward probabilities for actions were given according to

$$p(a) = 1 - [1 - p_0(a)]^{n(a)+1}$$

The value of $p_0(a)$, the baseline probability of reward for action a , was sampled from a uniform Bernoulli distribution and held fixed for the entire episode. $n(a)$ is the number of trials since action a was chosen. The sum of the baseline probabilities for both arms was always constant and equal to 0.5, so that $p_0(a_l) = 0.5 - p_0(a_r)$. The number of trials per episode varied uniformly from 50 to 100 (chosen randomly at the beginning of the episode). For training only, we held out the subset of environmental parameters $p_0 = 0.1-0.2$ and $0.3-0.4$, so the agent was only trained on $p_0 = [0.0-0.1] \cup [0.2-0.3] \cup [0.4-0.5]$. For this task, we used discount factor = 0.75 and trained for 3×10^6 steps ($\sim 15,000$ episodes).

After training, we tested the agent for 500 episodes with network weights held fixed. Analyses of model behavior were based on the original study, which in turn employed techniques introduced by Lau and Glimcher³³. We determined the influence of past choices and reward outcomes on trial-by-trial responses by fitting to a logistic regression model, as done by Lau and Glimcher³³ and Tsutsui et al.⁷, using a maximum trial lag of $N = 15$. Only trials from the last two-thirds of each episode were considered, in order to restrict our analyses to steady-state behavior.

To assess the amount of information encoded about each variable in the recurrent network, we calculated Spearman’s partial rank correlation between hidden unit activations and four factors: the last action, last reward, action \times reward interaction, and choice value. (Note that our task implementation abstracted over the distinction between actions and target objects, and we do not focus on the distinction between action and object value.) Because the values of the two actions were strongly anticorrelated ($r = -0.95$, Pearson’s correlation), we only considered the effect of the left action.

Simulation 2: adaptation of prefrontal-based learning to the task environment.

On each step the agent chose between two actions, a_l and a_r , and received a reward of R or 0. The magnitudes of R that could be obtained varied randomly across trials and independently between the two actions (distributed as $100 \times \beta$ (4.4)), and were cued to the agent on each trial. The reward probabilities of the two arms were perfectly anticorrelated and thus described by a single parameter $r(t)$, which evolved over time. In testing, the evolution of $r(t)$ mirrored the task

human subjects performed in Behrens et al.³⁴. There were two types of episodes. In episodes with low volatility first, $r(t)$ remained stable at 0.25 or 0.75 for 100 trials, and then alternated between 0.2 and 0.8 every 25 trials. In episodes with high volatility first, $r(t)$ alternated between 0.2 and 0.8 every 25 trials for 100 trials, and then remained stable at 0.25 or 0.75.

For training, we adapted the two-level generative model used for inference by Behrens et al.³⁴. Specifically, in training, $r(t)$ evolved according to $v(t)$, such that with probability $v(t)$, $r(t+1)$ was set to $1 - r(t)$, and with probability $1 - v(t)$, $r(t+1)$ remained the same as $r(t)$. Analogously, the volatility $v(t)$ also evolved over time. With probability k , $v(t+1)$ was set to $0.2 - v(t)$, and with probability $1 - k$, $v(t+1)$ remained the same as $v(t)$, where $\log(k)$ was sampled uniformly at the end of each episode between -4.5 and -3.5 . At the beginning of each training episode, $v(1)$ was randomly initialized to 0 or 0.2. When $v(t)$ was 0, $r(t)$ was set to 0.25 or 0.75, and when $v(t)$ was 0.2, $r(t)$ was set to 0 or 1. Tying $r(t)$ to $v(t)$ in this way matched the structure of the task subjects performed in Behrens et al.³⁴ (see above and Fig. 3a). The probabilistic model employed in that study did not impose this connection, with the result that dynamic adjustments in learning rate had less effect on expected reward. For completeness, we also ran simulations in which meta-RL was trained using data produced from the generative model employed by Behrens et al.³⁴ and obtained qualitatively similar results (data not shown).

Meta-RL was trained for 40,000 episodes and tested for 400 episodes, consisting of 200 trials each. Because of the length of these episodes, we did not include a fixation step. Training was conducted with learning rate = 5×10^{-5} , discount = 0.1, and an entropy cost linearly decreasing from 1 to 0 over the course of training. As in all simulations, weights were held fixed in testing.

For purposes of analysis, following Behrens et al.³⁴, we implemented a Bayesian agent, which performed optimal probabilistic inference under the generative model that defined the training experience. The Bayesian agent used a discretized representation of the posterior as in Behrens et al.³⁴. The 'estimated volatility' signal we report was the agent's expectation of volatility at time t . We also fit a multiple regression model to predict the Bayesian agent's estimated volatility at each time step from the recurrent network's hidden activations and applied this model to hidden activations on held-out episodes (in cross-validation) to generate the 'decoded volatility' signal we report for meta-RL.

Learning rates were also estimated for both LSTM and Bayes, based only on their behavior. To estimate learning rates, we fit a Rescorla–Wagner model to the behavior by maximum likelihood. Episodes of the two types (high volatility first and low volatility first) were fit separately. The Rescorla–Wagner model was fit to behavior from all episodes, but from only ten consecutive trials, sliding the window of ten trials in single-trial increments to observe how learning rates might change dynamically. Although only ten consecutive trials contributed directly to the likelihood, all trials previous to these ten nevertheless contributed indirectly through their influence on the evolving value estimate. As a control, we estimated learning rates for behavior generated by a Rescorla–Wagner model with a fixed learning rate, to ensure we could recover the true learning rate without bias.

To determine the proportion of hidden units coding for volatility, in each episode we regressed the activity of each hidden unit against the true volatility, coded as 0 or 1. (This regression model also included several covariates: previous reward, next action, and true probability of reward on the present time step. However, including these covariates had little effect on the proportion of units discovered as coding volatility.) We then counted units for which the slope of the volatility regressor was significantly different from zero, at a threshold of $P=0.05$, Bonferroni corrected for number of hidden units. The mean and standard error of this number are reported in the main text.

In a similar analysis, we also regressed across episodes at each time-step. In Supplementary Fig. 5a we plot over time the proportion of units for which the slope of the volatility regressor was significantly different from zero, at a Bonferroni corrected threshold of $P=0.05$. In Supplementary Fig. 5b we continuously vary the threshold across a wide range and plot the proportion of units, averaged over trials, exceeding the threshold.

Finally, we note that Behrens et al.³⁴ used the estimated volatility signal from the Bayesian model as a regressor, while we use the true volatility. We found that these two approaches did not produce qualitatively different results. However, it was slightly more difficult to interpret the Bayes estimated volatility signal at a fine-grained scale of individual trials because the Bayes estimated signal was sensitive to individual trial events. Particularly, early in an episode, the Bayes estimated signal was more a reflection of the specific outcomes that had been received than a reflection of volatility.

Simulation 3: reward prediction errors reflecting inferred value. Following Bromberg-Martin et al.¹², we trained our agent on a reversal task, in which the values of the two stimulus cues were anticorrelated, with one being rewarded and the other not rewarded. The rewarded stimulus switched intermittently in an uncued manner, such that there was a 50% chance of a switch at the beginning of each episode. Concretely, after the central fixation cue on step 1, the agent was presented with one of the two stimulus cues (tabular one-hot vectors) on step 2, indicating that it must either produce an action left (a_L) or an action right (a_R) on step 3. The reward was then delivered on step 4, followed by the start of the next trial. Failure to produce a valid response on any step gave a reward of -1 , while a

valid action resulted in a reward of 1 for the rewarded cue or 0 for the unrewarded cue and central fixation cue.

The rewarded and unrewarded cues were randomly determined at the beginning of the episode and held fixed for the duration of the episode, lasting 5 trials (20 steps). After training for 2×10^6 steps, or 1×10^5 episodes, using a backpropagation window of 200 steps, we tested for 1,200 episodes and calculated the reward prediction errors in response to stimulus presentation (step 2) for trials 1 and 2:

$$\delta_t = r_t + \gamma V_{t+1} - V_t$$

Baselines (i.e., respective V_t) were also extracted from steps 1 and 3 for comparison. We analyzed all test episodes in which a reward reversal had occurred with respect to the previous episode.

Simulation 4: 'model-based' behavior—the two-step task. The structure of the two-step task (Fig. 5a) was based directly on the version used by Miller et al.³⁹. After central fixation, on the first stage of the trial the agent chose either a_L or a_R and transitioned into one of two second-stage states S_1 or S_2 with probabilities $p(S_1|a_L) = p(S_2|a_R) = 0.8$ ('common' transition) and $p(S_1|a_R) = p(S_2|a_L) = 0.2$ ('uncommon' transition). The transition probabilities p were fixed across episodes. The states S_1 and S_2 yielded probabilistic reward according to $[p(r|S_1), p(r|S_2)] = [0.9, 0.1]$ or $[0.1, 0.9]$, with the specific reward contingencies having a 2.5% chance of randomly switching at the beginning of each trial. The agent was trained for 10,000 episodes of 100 trials each and tested with weights fixed on 300 further episodes. Behavioral analyses (Fig. 5c,e) closely followed corresponding procedures specified by Miller et al.³⁹ and Daw et al.¹⁵.

To compare the RPEs of our meta-RL agent to RPEs produced by model-free and model-based algorithms, we followed the procedure employed by Daw et al.¹⁵. As in that study, the reference model-free algorithm was SARSA(0) and the model-based algorithm was a prospective planner, both as implemented by Daw et al.¹⁵. However, instead of performing Q-learning at the second stage, which is not optimal, the model-based algorithm tracked the latent variable H describing whether $p(r|S_1)$ was 0.9 (equivalent to whether $p(r|S_2)$ was 0.1). The algorithm updated its belief about H following each outcome, using knowledge of the task's true structure. Both algorithms were applied to the same trial sequence experienced by the meta-RL agent. Each agent generated a reward prediction error on each step, as described by Daw et al.¹⁵, but with the expectation of reward at the second stage being $p(H) \times 0.9 + (1 - p(H)) \times 0.1$ if the agent was in state S_1 and $(1 - p(H)) \times 0.9 + p(H) \times 0.1$ if the agent was in state S_2 .

RPEs from meta-RL were calculated as in simulation 3. We were interested in determining whether meta-RL's RPEs were model-based only, model-free only, or a mixture of both. This analysis must take account of the fact that model-free and model-based RPEs themselves are somewhat correlated (Supplementary Fig. 6). We first regressed meta-RL's RPEs directly against the model-based RPEs of the prospective planner. Finding an almost perfect correlation, we then asked whether meta-RL's RPEs were at all related to the component of model-free RPEs that is orthogonal to the model-based RPEs. We did this by regressing the model-based RPE out of the model-free RPE and looking for a correlation of meta-RL's RPE with the residual. We note that, unlike subjects in the fMRI study by Daw et al.¹⁵, our model showed a purely model-based pattern of behavior (see Fig. 5c,e). Correspondingly, we found exclusively model-based RPEs in the present analysis. As a confirmatory analysis, we also performed a multiple regression with meta-RL's RPEs as the dependent variable, and both model-based and model-free RPEs as independent variables (Supplementary Fig. 6).

Simulation 5: learning to learn. An earlier version of this simulation was presented by Wang et al.³¹. To simulate Harlow's⁴⁰ learning-to-learn task, we trained our agent with RGB pixel input using the Psychlab framework⁵⁴. An 84×84 pixel input represented a simulated computer screen (see Fig. 6a and Supplementary Video 1). At the beginning of each trial, this display was blank except for a small central fixation cross. The agent selected discrete left–right actions that shifted its view approximately 4.4 degrees in the corresponding direction, with a small momentum effect (alternatively, a no-op action could be selected). The completion of a trial required performing two tasks: saccading to the central fixation cross, followed by saccading to the correct image. If the agent held the fixation cross in the center of the field of view (within a tolerance of 3.5 degrees visual angle) for a minimum of four time steps, it received a reward of 0.2. The fixation cross then disappeared and two images, drawn randomly from the ImageNet dataset⁵⁵ and resized to 34×34 , appeared on the left and right side of the display, respectively (Fig. 6a). The agent's task was then to 'select' one of the images by rotating until the center of the image aligned with the center of the visual field of view (within a tolerance of 7 degrees visual angle). Once one of the images was selected, both images disappeared and, after an intertrial interval of 10 time-steps, the fixation cross reappeared, initiating the next trial. Each episode contained a maximum of six trials or 3,600 steps. Each selected action was mandatorily repeated for a total of four time-steps (as in Mnih et al.⁵²), meaning that selecting an image took a minimum of three independent decisions (twelve primitive actions) after having completed the fixation. It should be noted, however, that the rotational position of the agent was not limited; that

is, 360 degree rotations could occur, while the simulated computer screen only subtended 65 degrees.

Although new ImageNet images were chosen at the beginning of each episode (sampled with replacement from a set of 1,000 images), the same two images were reused across all trials within an episode, though in randomly varying left–right placement, similarly to the objects in Harlow's⁴⁰ experiment. And as in that experiment, one image was arbitrarily chosen to be the 'rewarded' image throughout the episode. Selection of this image yielded a reward of 1.0 while the other image yielded a reward of -1.0. The LSTM was trained using backpropagation through time with 100-step unrolls.

During test, network weights were held fixed (in convolutional network and LSTM-A3C) and ImageNet images were drawn from a separate, held-out set of 1,000 never presented during training. Learning curves were calculated for each replica by taking a rolling average of returns over episodes and normalizing from 0% to 100%, where 50% represented performance of randomly choosing between the two images and 100% represented optimal performance. We trained 50 networks using hyperparameters adopted from Wang et al.³¹ (learning rate = 0.00075, discount = 0.91, entropy cost $\beta_1 = 0.001$, and state-value estimate cost $\beta_2 = 0.4$) and found that 28% reached maximal performance after $\sim 1 \times 10^5$ episodes (per thread, 32 threads). Because this task required not only learning abstract rule structure, but also processing visually complex images, staying oriented toward the simulated computer screen, and central fixation, much of the initial training time was spent at 0% performance. Ultimately successful agents varied in the number of episodes required to overcome these initial hurdles (median 3,413, range 1,810–29,080). To characterize acquisition of abstract task structure dissociated from these other considerations, we identified two change points in the learning curves (achieving above-chance level performance and approaching ceiling performance) and plotted average reward as a function of trial number within the episode at various quantiles of performance interpolated between these change-points (six quantiles) as well as the final performance at end of training. Agents also varied in the number of episodes required to achieve maximal performance after surpassing chance-level performance (median 4,875, range 1,500–82,760).

To be sure that the meta-RL network was truly performing the role-filler assignment that is intended to be required by the Harlow⁴⁰ task, rather than relying on a simpler strategy, we trained the same model on a variant of the Harlow⁴⁰ task wherein the non-target (unrewarded) image changed across trials, always taking the form of an unfamiliar picture. Following training, the performance of our network was identical to its performance in the original task, with chance accuracy on trial 1 followed by near ceiling accuracy for the remainder of the same trial block (data not shown).

Simulation 6: the role of dopamine—effects of optogenetic manipulation. We implemented a probabilistic risk/reward task as described in Stopper et al.⁴⁴, in which subjects chose between a 'safe' arm that always offered a small reward ($r_S = 1$) or a 'risky' arm that offered a large reward ($r_L = 4$) with a probability $p(a_{\text{large}}) = 0.125$ (Safe Arm Better block) or 0.5 (Risky Arm Better block), sampled at the beginning of the episode. In direct analogy with the original study, the agent was first required to make 5 forced pulls each of the safe and risky arms (in randomized pairs), followed by 20 free pulls. The locations (i.e., associated actions) of the risky and safe arms could switch from episode to episode, with a probability of 0.5.

To simulate optogenetic stimulation, it was necessary to implement a variant of our original architecture in which two separate LSTMs (48 units each) model value estimate and policy (rather than a single LSTM with two linear outputs), in alignment with standard actor/critic architectures associated with basal ganglia and PFC function (see Supplementary Fig. 3). Concretely, the critic is modeled as an LSTM that takes as input the state observation, last reward received and last action taken, outputting value estimate. The actor receives the state observation, the last action taken and the reward prediction error that is computed on the basis of the

value estimate output from the critic, and subsequently outputs the policy (see Supplementary Fig. 9). Because the RPE is computed on the basis of the prefrontal network value output, it could not be fed as an input to the critic, given our other architectural and algorithmic assumptions.

Optogenetic stimulation was simulated by manipulating the value of the reward prediction error fed into the actor, in directly analogy with lateral habenula and ventral tegmental area stimulation in Stopper et al.⁴⁴. After training our agent in the control condition in which normal RPEs were input to the actor (4×10^6 steps; 6.67×10^4 episodes of 30 trials each), we tested our agent for 1,333 episodes in four different conditions designed to approximate the different stimulation protocols employed by Stopper et al.⁴⁴: (1) control (analogous to Stopper's baseline no-stimulation condition), (2) block risky reward—if the risky arm was chosen and rewarded, subtract 4 from the RPE (approximating lateral habenula stimulation with 4 trains of pulses), (3) block safe reward—if the safe arm was chosen, subtract 1 from the RPE (approximating lateral habenula stimulation with 1 train of pulses), and (4) block risky loss—if the risky arm was chosen and not rewarded, add a reward of 1 to the RPE (approximating ventral tegmental area stimulation).

Statistics. Where we report variance in our simulations that is parallel to variance in corresponding human or animal experiments, we ran $n = 8$ replicas of the simulation unless otherwise stated. All replicas were identical except for having a different random seeds. We report variance over these replicas with the idea that this is most comparable to animal experiments. Each replica also contained many episodes, which were used to perform detailed analyses in order to compare to human and animal data. Where we did this, other runs produced very similar patterns of results.

We report mean and s.e.m. throughout the paper unless otherwise noted. Correlation was by Pearson's r unless otherwise noted. All tests were two-tailed unless otherwise noted. The data distribution was assumed to be normal, but this was not formally tested. No data points were excluded from the analyses unless otherwise noted. No statistical methods were used to predetermine sample sizes, since all networks were able to be trained to full performance. Data collection and analysis were not performed blind to the conditions of the experiments, as this did not apply to our simulations. Data collection and assignment to experimental groups also did not apply, since all networks were equivalent before training.

Reporting Summary. Further information on experimental design is available in the Nature Research Reporting Summary.

Code and data availability. Data sharing is not applicable to this study since no experimental datasets were generated or analyzed, but behavioral and activations data from trained versions of the models from simulations 1–6, together with open-sourced versions of analysis scripts, are available from the corresponding author upon request.

References

- Hochreiter, S. & Schmidhuber, J. Long short-term memory. *Neural Comput.* **9**, 1735–1780 (1997).
- Mnih, V. et al. Asynchronous methods for deep reinforcement learning. in *Proc. 33rd Intl. Conf. Machine Learning* **48**, 1928–1937 (JMLR, New York, 2016).
- Graves, A., Jaitly, N. & Mohamed, A.-r. Hybrid speech recognition with deep bidirectional LSTM. in *IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU) 2013* 273–278 (IEEE, 2013).
- Leibo, J. Z. et al. Psychlab: a psychology laboratory for deep reinforcement learning agents. Preprint at <https://arxiv.org/abs/1801.08116> (2018).
- Deng, J. et al. ImageNet: a large-scale hierarchical image database. in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2009* 248–255 (IEEE, 2009).

Reporting Summary

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

Statistical parameters

When statistical analyses are reported, confirm that the following items are present in the relevant location (e.g. figure legend, table legend, main text, or Methods section).

n/a Confirmed

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

Our web collection on [statistics for biologists](#) may be useful.

Software and code

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

Data collection

TensorFlow code developed by the authors was used in all simulations, as fully described in the paper.

Data analysis

Matlab and python code written by the authors was used for all analyses, as fully described in the paper.

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

Data

Policy information about [availability of data](#)

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

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

Our paper presents no experimental data. Our code availability statement indicates how interested readers can replicate our analyses.

Field-specific reporting

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

Life sciences Behavioural & social sciences

For a reference copy of the document with all sections, see [nature.com/authors/policies/ReportingSummary-flat.pdf](https://www.nature.com/authors/policies/ReportingSummary-flat.pdf)

Behavioural & social sciences

Study design

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

Study description	No experimental work is reported. Only computer simulations are presented.
Research sample	There was no research sample. We only report computer simulations.
Sampling strategy	This is not relevant. We only report computer simulations.
Data collection	This is not relevant. We only report computer simulations.
Timing	This is not relevant. We only report computer simulations.
Data exclusions	This is not relevant. No data was excluded.
Non-participation	This is not relevant. There were no participants.
Randomization	This is not relevant. There was nothing to randomize.