

1 **Behavioural and computational evidence for memory**
2 **consolidation biased by reward-prediction errors**

3
4 **Author names:** Emma L. Roscow*¹; Matthew W. Jones¹; Nathan F. Lepora²

5 **Author affiliations:** ¹ School of Physiology, Pharmacology & Neuroscience, University of Bristol, University
6 Walk, Bristol BS8 1TD, UK; ² Department of Engineering Mathematics, University of Bristol, Bristol, UK

7 * corresponding author

8 **Abstract**

9 Neural activity encoding recent experiences is replayed during sleep and rest to promote consolidation of the
10 corresponding memories. However, precisely which features of experience influence replay prioritisation to
11 optimise adaptive behaviour remains unclear. Here, we trained adult male rats on a novel maze-based rein-
12 forcement learning task designed to dissociate reward outcomes from reward-prediction errors. Four variations
13 of a reinforcement learning model were fitted to the rats' behaviour over multiple days. Behaviour was best
14 predicted by a model incorporating replay biased by reward-prediction error, compared to the same model with
15 no replay; random replay or reward-biased replay produced poorer predictions of behaviour. This insight dis-
16 entangles the influences of salience on replay, suggesting that reinforcement learning is tuned by post-learning
17 replay biased by reward-prediction error, not by reward per se. This work therefore provides a behavioural and
18 theoretical toolkit with which to measure and interpret replay in striatal, hippocampal and neocortical circuits.

19 **1 Introduction**

20 To make good decisions, it is typically beneficial to use past experience to guide future behaviour. Actions
21 which have previously produced good outcomes in a similar context can be reinforced to adapt behaviour for
22 maximising benefit. Crucial to this mechanism is the ability for neuronal spiking activity to drive synaptic plas-
23 ticity, strengthening the synaptic connections between neurons to establish functional networks which encode
24 task-relevant information or drive task-relevant actions. These functional networks are refined during sleep
25 and rest, when many neurons switch to an “offline” state in which they replay activity encoding previous or
26 anticipated upcoming experiences rather than current events or behaviours (Yu et al. 2017). This offline replay,
27 found across cortical, limbic and basal ganglia regions, has been suggested to play a role in decision-making
28 (Pfeiffer and Foster 2013), emotional processing (Cairney et al. 2014), generalising across episodes (Lewis and
29 Durrant 2011), and reinforcement learning (Dupret et al. 2010).

30 Studies in which replay has been manipulated provide strong evidence for its contributions to memory consol-
31 idation. Artificially enhancing replay by presenting odours or sounds during sleep, which had previously been
32 paired with object locations or visual stimuli, leads to better subsequent recall of the paired stimuli (Rasch et al.
33 2007; Rudoy et al. 2009; Antony et al. 2012; Bendor and Wilson 2012). Disrupting replay events, meanwhile,
34 impairs subsequent spatial memory (Girardeau et al. 2009; Ego-Stengel and Wilson 2010; Jadhav et al. 2012;
35 Michon et al. 2019).

36 An examination of how replay aids these cognitive processes requires assessment of which activity is replayed
37 with greatest strength or frequency. Activity which is associated with experiences of reward (Foster and Wilson
38 2006; Lansink et al. 2009; Singer and Frank 2009) or fear (Girardeau et al. 2017; Wu et al. 2017), or with
39 recent experiences (Cheng and Frank 2008), is replayed preferentially. This suggests a replay bias towards the
40 most salient experiences to be processed, consolidated or incorporated into the internal model of the world.
41 However, these salient experiences could also be interpreted as those with the highest prediction error, i.e.
42 the most informative experiences for updating internal models and for reinforcement learning. Tasks which
43 involve learning the locations of rewards often conflate reward with reward-prediction error (RPE), leading to
44 the possibility that apparent replay biases towards reward actually reflect biases towards RPE.

45 Here we explore the possibility that it is reward prediction errors, rather than reward or salience, which biases
46 replay. We used variations of a reinforcement learning model, Q-learning, to estimate the value of actions
47 encoded in the striatum during a reinforcement learning task, and varied the amount and type of replay in the
48 model to predict behaviour. In the striatum, representations of reward values differ following learning acquired
49 over weeks compared to when acquired over minutes (Wimmer et al. 2018), and, correspondingly, reward-
50 responsive cells are replayed preferentially in the ventral striatum (Lansink et al. 2009). We therefore propose
51 that replay triggers value updates in the striatum, to enhance striatum-dependent reinforcement learning, and
52 moreover that activity encoding events that resulted in high RPE is preferentially replayed.

53 Q-learning (Watkins 1989) has been used successfully to model reinforcement learning, particularly in humans
54 (O’Doherty et al. 2003, Daw et al. 2005) but also in rodents (Kim et al. 2013, Ito and Doya 2009). Q-learning
55 models fit both behavioural outcomes and striatal activity, suggesting that they describe mechanisms of updating
56 values in the striatum in response to RPEs which in turn guide behaviour (Day et al. 2014, Morris et al. 2010,

57 Pagnoni et al. 2002, Roesch et al. 2007). Temporal-difference-based RPEs, i.e. the difference between expected
58 reward and actual reward which drives the update of Q-values, resemble quite closely the dopaminergic input
59 of ventral tegmental area (VTA) to the striatum (McClure et al. 2003, Roesch et al. 2007, Schultz 2016),
60 which mediates synaptic plasticity in the striatum (Calabresi et al. 2007) and may provide a mechanism of
61 biological implementation of Q-learning. Dyna-Q (Sutton 2014), a variant of Q-learning which incorporates
62 offline temporal-difference updates, has been used to model replay in ways which produce learning qualitatively
63 similar to animal reinforcement learning (Johnson and Redish 2005). RPE-biased replay incorporated into
64 machine learning algorithms show that it can also be very efficient, learning to play Atari games (Andrychowicz
65 et al. 2017) or navigate a simulated environment (Karimpanal and Bouffanais 2017) faster and with more
66 success compared to replay without such a bias.

67 We trained 6 rats on a stochastic reinforcement learning task which elicited both positive and negative RPE, and
68 fitted Q-learning parameters to each rat's behavioural data. We then included replay events between sessions,
69 to simulate the effect of replay during sleep on reinforcement learning. Four replay policies were compared,
70 prioritising state-action pairs to be updated according to different biases: random replay, replay proportional to
71 expected reward, and two forms of RPE-biased replay. Random replay was included as a control, while reward-
72 biased replay reflects the prevailing view of how replay is prioritised. Fitting the model parameters showed that
73 the two RPE-biased replay policies increased the model's predictive accuracy, while random and reward-biased
74 replay impaired model performance. This suggests that replay between sessions of a probabilistic reinforcement
75 learning task in rats is biased by RPE and not by reward.

76 **2 Results**

77 **Animals successfully learned a stochastic reinforcement learning task**

78 Six rats were trained to forage for stochastic sucrose rewards on a three-armed maze, to assess their reinforce-
79 ment learning on a task where reward outcome and reward-prediction error (RPE) were dissociable. Each arm
80 was assigned as either "high probability", "mid probability" or "low probability", which determined the pro-
81 tocol for reward delivery (fig. 1a). For the first 15 training sessions, the high-probability arm delivered a reward
82 on 75% legitimate entries to the arm, the mid-probability arm on 50%, and the low-probability arm on 25%. A
83 legitimate entry was one in which a different arm had been entered on the previous trial; entering the same arm
84 twice in a row was incorrect and did not result in a reward delivery. For sessions 16-20, the reward probabilities
85 for the high- and low-probability arms were amplified: reward was delivered on 87.5% and 12.5% legitimate
86 entries respectively. For sessions 21-22 the reward probabilities for the high- and low-probability arms were
87 switched, such that the (formerly) high- and low- probability arms delivered reward on 12.5% and 87.5% of
88 legitimate entries respectively. This set-up meant that receiving a reward in a low-probability arm would elicit
89 a higher RPE than the same reward value in a high-probability arm, so reward outcome and RPE could be
90 dissociated.

91

92 Over 22 sessions, animals learned to distinguish between the high-, mid- and low-probability arms in their
93 frequency of visits to each arm, indicating successful learning of the reward probabilities. Rats performed 45.1
94 ± 2.5 trials per session, eventually showing a significant preference for the high-probability arm and against
95 the low-probability arm, evident by session 6 and stable by session 10. The six animals varied in the degree of
96 their discrimination between the arms (fig. 1b), but on average they distinguished between all arms on 13 out
97 of 22 sessions (fig. 1c; χ^2 test, Bonferroni-corrected), visiting the arms which delivered a higher probability of
98 reward more often, primarily in later sessions. The differences in arm discrimination between animals may be
99 accounted for by the orientation of the maze in the room; for example, animals may have shown a confounding
100 preference for the arm which was closest to the door of the recording room, an effect which was overcome by
101 rotating the arm probabilities between animals. (For rats H and K, the mid-probability arm was closest to the
102 door; for rats I and L, the high-probability arm was closest to the door; and for rats J and M, the low-probability
103 arm was closest to the door.)

104 To quantify performance on the task, each trial was coded as optimal or suboptimal according to the animal's
105 choice of arm given the arm most recently visited. Because no reward was given for re-entering the same arm
106 visited on the previous trial, the optimal action choice following a visit to the mid- or low-probability arm was
107 to visit the high-probability arm; the optimal action following the high-probability arm was the mid-probability
108 arm. Over sessions, animals increased the proportion of trials on which they behaved optimally, achieving
109 performance significantly above chance from session 3 onwards (fig. 1d, 46 trials optimal out of 106, $p = 0.02$,
110 binomial test, Bonferroni-corrected).

111 Reward probabilities were changed twice over the course of learning, triggering clear changes in behaviour. In

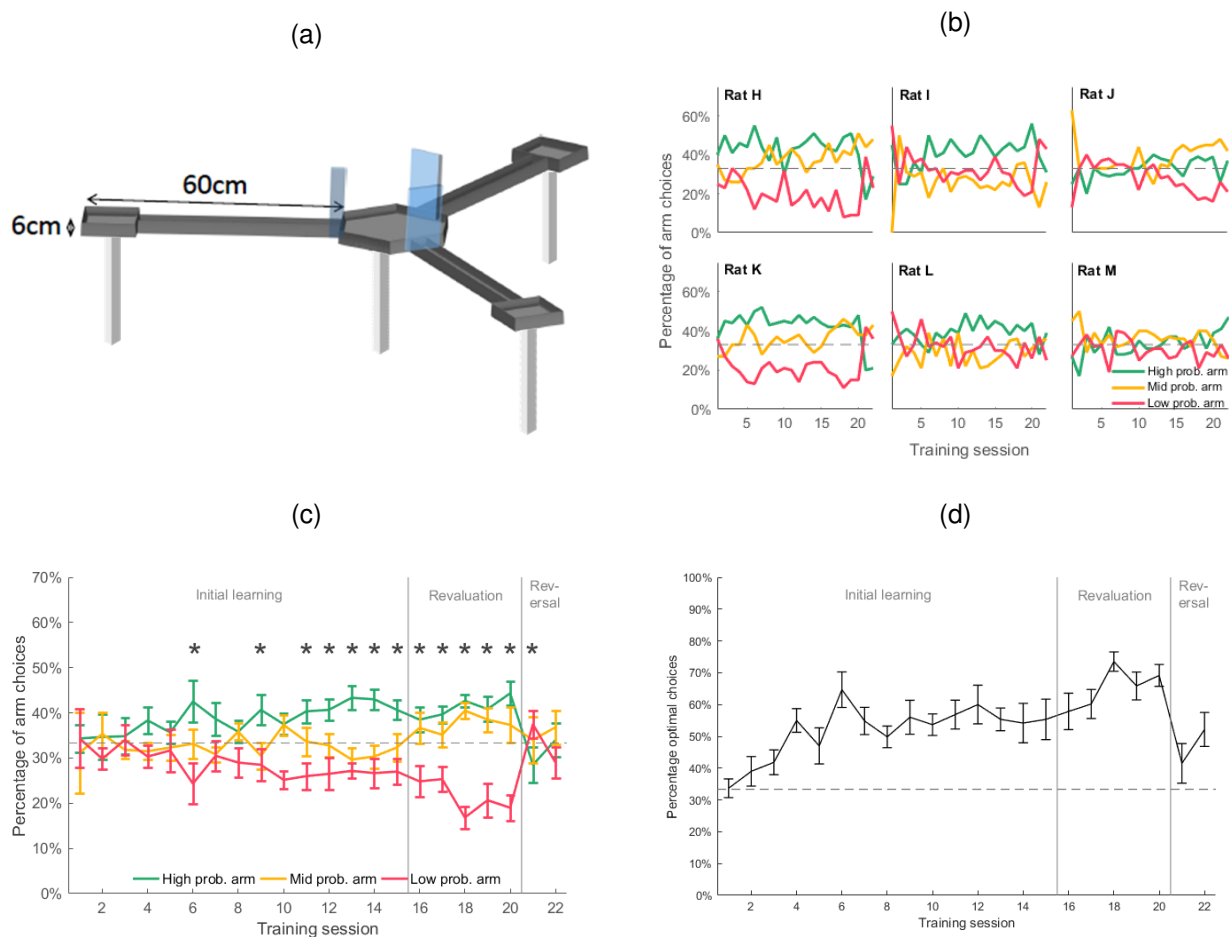


Figure 1: **A.** Illustration of the maze used to train animals. Lick ports located at the end of each arm delivered reward with either high, medium or low probabilities. **B.** Frequency of entry to each arm over all sessions, shown separately for each rat. **C.** Average frequency of entry to each arm. * indicates arm choices statistically different from each other (χ^2 test). **D.** Mean proportion of trials on which the optimal arm was chosen, according to highest probability of reward. Dashed lines represent chance level (33.3%). Error bars represent standard error of the mean (s.e.m.).

112 the revaluation learning stage (sessions 16-20), the reward probabilities at each arm became more distinct: the
 113 high-probability arm delivering an 87.5% probability of reward compared to 75% in the initial learning stage,
 114 and the low-probability arm delivering a 12.5% probability of reward compared to 25% in the initial learning
 115 stage. This change offered a higher incentive to avoid the low-probability arm and, correspondingly, preference
 116 for the high-probability arm over the low-probability arm increased compared to the previous five sessions
 117 (fig. 1c; repeated-measures ANOVA, $F = 8.7$, $p = 0.006$). As a result, the rate of optimal performance was also
 118 greater in the revaluation stage than the last five sessions of the initial learning stage (fig. 1d; repeated-measures
 119 ANOVA, $F = 15.2$, $p = 0.001$).

120 The definition of optimal behaviour was the same in the initial and revaluation learning stages, because the
 121 arms did not change. However, optimal behaviour required a different behavioural policy in the reversal learn-
 122 ing stage (sessions 21-22) when the high- and low-probability arms were switched. As expected, optimal
 123 performance correspondingly dipped when reward probabilities were reversed in sessions 21 to 22 as this new
 124 behavioural policy was learned. The frequency of optimal arm choices was lower for the reversal learning stage
 125 than the last two sessions of the revaluation stage (repeated-measures ANOVA post-hoc test, $p = 0.17$), although
 126 it did not differ significantly from the last two sessions of the initial learning stage (repeated-measures ANOVA

127 post-hoc test, $p > 0.05$). These behavioural data demonstrate that reward probabilities successfully influenced
128 learning and behaviour in the task, and that animals were capable of showing flexibility in response to chan-
129 ging reward. We therefore went on to test whether reinforcement learning algorithms were able to recapitulate
130 rat behaviour and whether instantiating between-session ("offline") replay of different task features improved
131 model performance.

132 Q-learning modelled animal behaviour

133 We trained a Q-learning algorithm with no replay to generate probabilities of each action for each trial, based
134 on Q-values estimated from the animals' previous experience (fig. 2). Q-learning is a reinforcement learning
135 algorithm in which an agent selects actions in its environment and observes the outcome, recording at each time
136 step t its starting state s_t , selected action a_t , resulting reward r_t , and resulting state s_{t+1} . The agent builds up a
137 matrix Q of Q-value estimates for every state-action pair:

$$\begin{bmatrix} Q_{s_1, a_1} & Q_{s_1, a_2} & \cdots & Q_{s_1, a_A} \\ Q_{s_2, a_1} & Q_{s_2, a_2} & \cdots & Q_{s_2, a_A} \\ \vdots & \vdots & \ddots & \vdots \\ Q_{s_S, a_1} & Q_{s_S, a_2} & \cdots & Q_{s_S, a_A} \end{bmatrix} \quad (1)$$

138 corresponding to the future discounted expected reward, i.e. the temporal difference between the current state
139 and the reward state. These Q-value estimates are used to guide actions to maximise reward. At each time step
140 t , the Q-value for the state-action pair observed is updated by:

$$Q(s_t, a_t) \leftarrow (1 - \alpha) \cdot Q(s_t, a_t) + \alpha \cdot (r_t + \gamma \cdot \max_a Q(s_{t+1}, a)) \quad (2)$$

141 where $\alpha \in (0, 1)$ is a learning rate parameter which determines the degree to which new information overrides
142 old information, and $\gamma \in (0, 1)$ is a discount parameter which determines the importance of long-term gains.

143 In this task, entries into a chosen arm (and arrival at the goal location at the end of the arm) were modelled
144 as actions, while the arm entered on the previous trial, on which reward probabilities were contingent, were
145 modelled as states. Each trial therefore gave rise to one state-action transition out of nine possible state-action
146 pairs.

147 For each trial, a matrix of Q-values for all state-action pairs was updated based on experience and used to
148 calculate predicted action probabilities, which were compared to the observed frequencies of state-action pairs
149 to produce a vector of errors for the three available actions. A reliability error was calculated from the summed
150 square of the error vector, weighted by the prevalence of the state. This produced a measure of how reliably the
151 Q-value estimates predicted behaviour (fig. 2; see Materials and Methods).

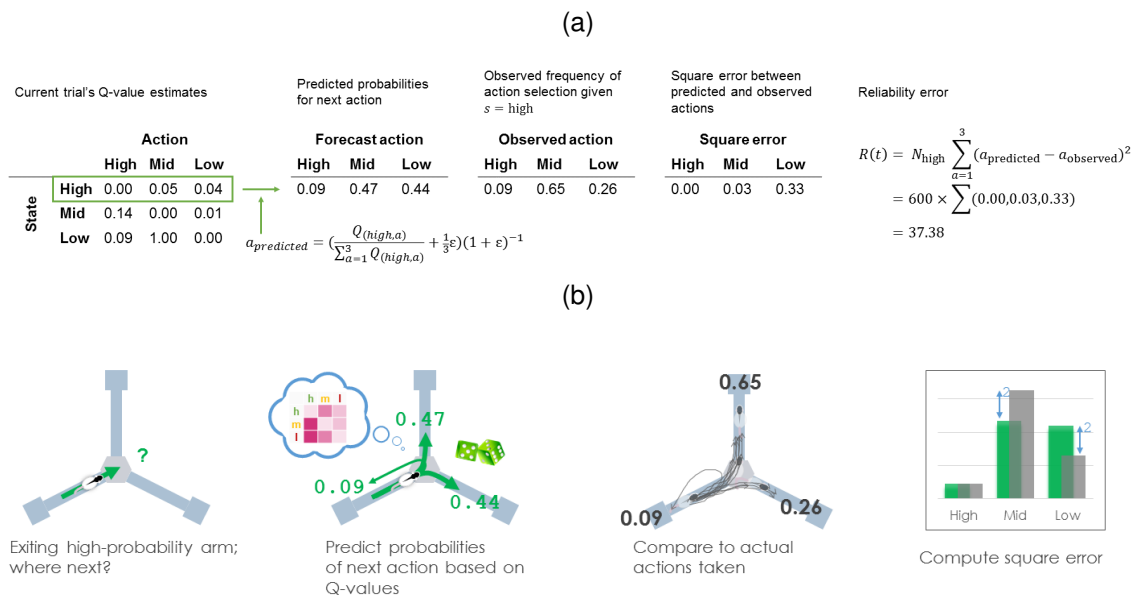


Figure 2: Example of model prediction for one trial, $t = 100$, in which rat H had most recently visited the high-probability arm ($s = \text{high}$) and chose the mid-probability arm ($a = \text{mid}$). **A.** The far left table shows the Q-learning model's estimate of the Q-values based on rat H's experience to date. Other tables show the predicted action probabilities calculated from the Q-values, the ground-truth of observed action frequencies over all visits to this state, and the mean square error between them. Far right shows how the error for this trial is calculated. **B.** A cartoon illustration of the same trial.

152

153 Observed action frequency correlated well with predicted action probabilities (fig. 3a), indicating a good
 154 baseline model for reinforcement learning. Predicted action probabilities were binned in 100 percentile-bins
 155 for each animal, and for each bin the average frequency of these actions occurring was compared to the average
 156 predicted probability, resulting in a strong correlation ($r = 0.92$, $p = 7.8e^{-08}$, Pearson's correlation). This result
 157 was consistent across animals (correlations ranging from $r = 0.86$ to $r = 0.96$).

158 The error between predicted action probability and observed action frequency spanned a large range, which
 159 was greatest in the earlier training sessions and diminished towards 0 for later training sessions as Q-values
 160 were learned (fig. 3b; early trials in blue have larger errors).

161

162 Reliability errors spanned a different range for each animal (fig. 3c), so all further analysis was performed
 163 on reliability errors normalised by the mean reliability error for each animal. On this measure, normalised
 164 reliability errors were similarly highest in early training sessions, when behaviour is least optimal and most
 165 unpredictable. Following this, reliability errors became consistently low for most sessions (fig. 3d), confirming
 166 a consistent fit with behaviour which captured the learning process over multiple sessions and changes in reward
 167 probabilities.

168 As described in Materials and Methods, the reliability error was used as the cost function to optimise three
 169 parameters in the Q-learning algorithm for each animal: a learning rate α , a discount factor γ , and an exploration
 170 factor ϵ . The resulting optimised parameter values are shown in table 1. A perturbation analysis was performed

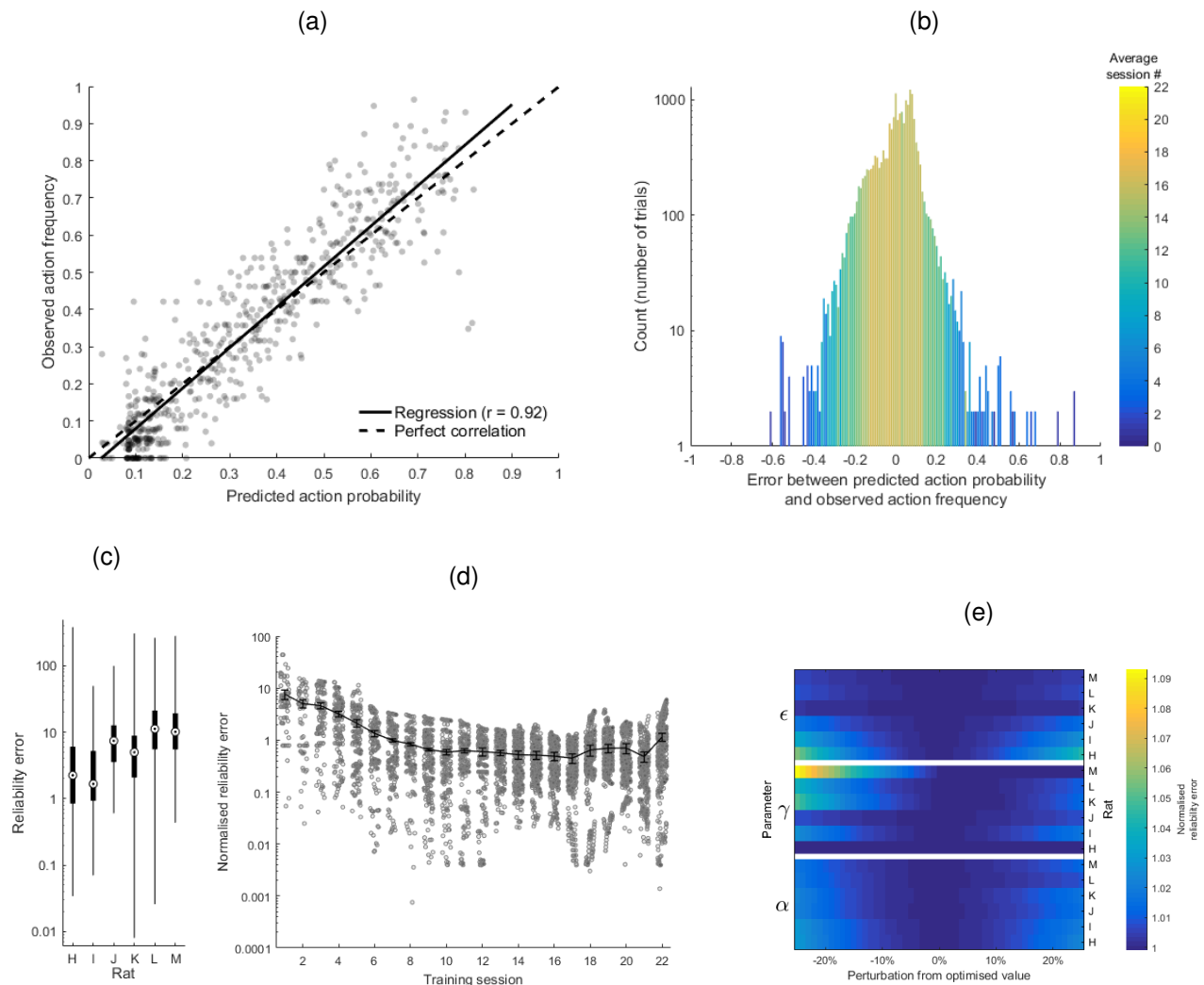


Figure 3: **A.** Reliability diagram (trials pooled across all animals). Observed frequency indicates how often an action was chosen by the animal, averaged over similar predicted action probabilities. Colour scale indicates the average session in which the predicted action probabilities occurred, for each bin; blue data points indicate predicted actions generally early in learning while yellow data points indicate predicted actions generally late in learning. Solid line represents regression ($r = 0.92$, $p=7.8e^{-244}$); dashed line indicates perfect correlation. **B.** Histogram of residuals of the data in A. Colour scale indicates on average what session the residuals within each bin occurred in. **C.** Range of reliability errors (calculated from residuals) for each animal. A reliability error of 0 reflects perfect modelling of action choices. Boxes represent 25th and 75th percentiles, circles represent median. **D.** Reliability errors for each trial grouped into training sessions, normalised to the average reliability error for each animal (shown in table 1). Data points show normalised reliability error for all trials; solid line represents mean for all animals. Error bars represent s.e.m. **E.** Change in reliability error, normalised to the optimised reliability error for each animal, with varying perturbations to the optimised parameter values. The optimised values for learning rate α , discount factor γ and exploration factor ϵ were individually perturbed by 1%-25% above and below the optimised value and the Q-learning algorithm was trained on behavioural data according to the perturbed parameter values 1,000 times to obtain an average.

171 to verify that the Q-learning results were sufficiently insensitive to perturbations to the optimised parameter
 172 values. At the optimised values, the average normalised reliability error over all trials was, by definition, 1.
 173 Perturbing these values by up to 25% in either direction increased the normalised reliability error by less than
 174 0.05 in most cases (fig. 3e) and less than 0.1 in all cases, indicating that reliability errors were not overly
 175 sensitive to small changes in parameter values.

	α	γ	ϵ	Reliability error
Rat H	0.009470	3.340e-09	0.3451	8.688
Rat I	0.01399	0.2972	0.4035	4.355
Rat J	0.02591	0.5153	0.3173	10.08
Rat K	0.06887	1.000	0.09363	10.66
Rat L	0.6522	1.000	0.3117	18.72
Rat M	0.1345	1.000	0.3137	16.92

Table 1: Optimised parameter values for Q-learning algorithm trained on each animal's behavioural data. α is the learning rate, γ is the discount factor, and ϵ is the exploration factor.

176 In summary, the Q-learning algorithm proved able to recapitulate rat behaviour over the course of training and
177 adaptation to new task conditions. The model was robust across a range of parameter values and established a
178 sound basis on which to quantify the effects of mimicking replay by updating Q values between sessions.

179 **Adding RPE-biased replay to the Q-learning model improved prediction** 180 **accuracy, whereas reward-biased and random replay both reduced ac-** 181 **curacy**

182 Against the baseline of no-replay, a variant of the Q-learning algorithm with replay was trained on the same
183 data, with a specified number of samples chosen from all the trials experienced so far to be replayed between
184 each session. Q-learning parameters were optimised for a fixed ($1 \leq n \leq 100$) number of replay events
185 between each session, for each replay policy. All trials experienced by the animal were stored in a memory
186 buffer, and for each replay event a state-action pair was chosen according to the replay policy and a sample
187 trial from this state-action pair was used to update its Q-value. With a random replay policy, all state-action
188 pairs that had been experienced were sampled at random. With a reward-biased replay policy, state-action pairs
189 were sampled in proportion to their Q-values, so that state-action pairs at which rewards had been experienced
190 most frequently would be replayed most. With an RPE-prioritised replay policy, the state-action pair with the
191 highest recent average RPE was sampled. With an RPE-proportional replay policy, state-action pairs were
192 sampled in proportion to their recent average RPE. These latter policies offered two variations on preferentially
193 updating state-action value(s) which had generated the greatest errors, concentrating efforts on correcting the
194 most erroneous expectations of reward.

195 Compared to the no-replay Q-learning baseline, replay biased by RPE produced a more reliable model of learn-
196 ing, while replay that was random or biased by reward produced a less reliable model (fig. 4a; orange and
197 purple compared to blue and green). Both the random and reward-biased replay policies resulted in higher reli-
198 ability errors ($p=8.8e^{-11}$ random, $p=1.6e^{-08}$ reward-biased, Wilcoxon signed rank test, Bonferroni-corrected),
199 even with a small amount of replay. Conversely, both the RPE-biased replay policies resulted in lower reli-
200 ability errors ($p = 6.6e^{-12}$ RPE-prioritised, $p = 6.3e^{-10}$ RPE-proportional). This was true even when one
201 additional sample was replayed between sessions (fig. 4b) and remained true when more samples were re-
202 played between sessions (fig. 4c-4e). Replay of information encoded during trials associated with the most
203 unexpected outcomes therefore significantly improved learning in the model, whereas replay of rewarded trials
204 proved detrimental.

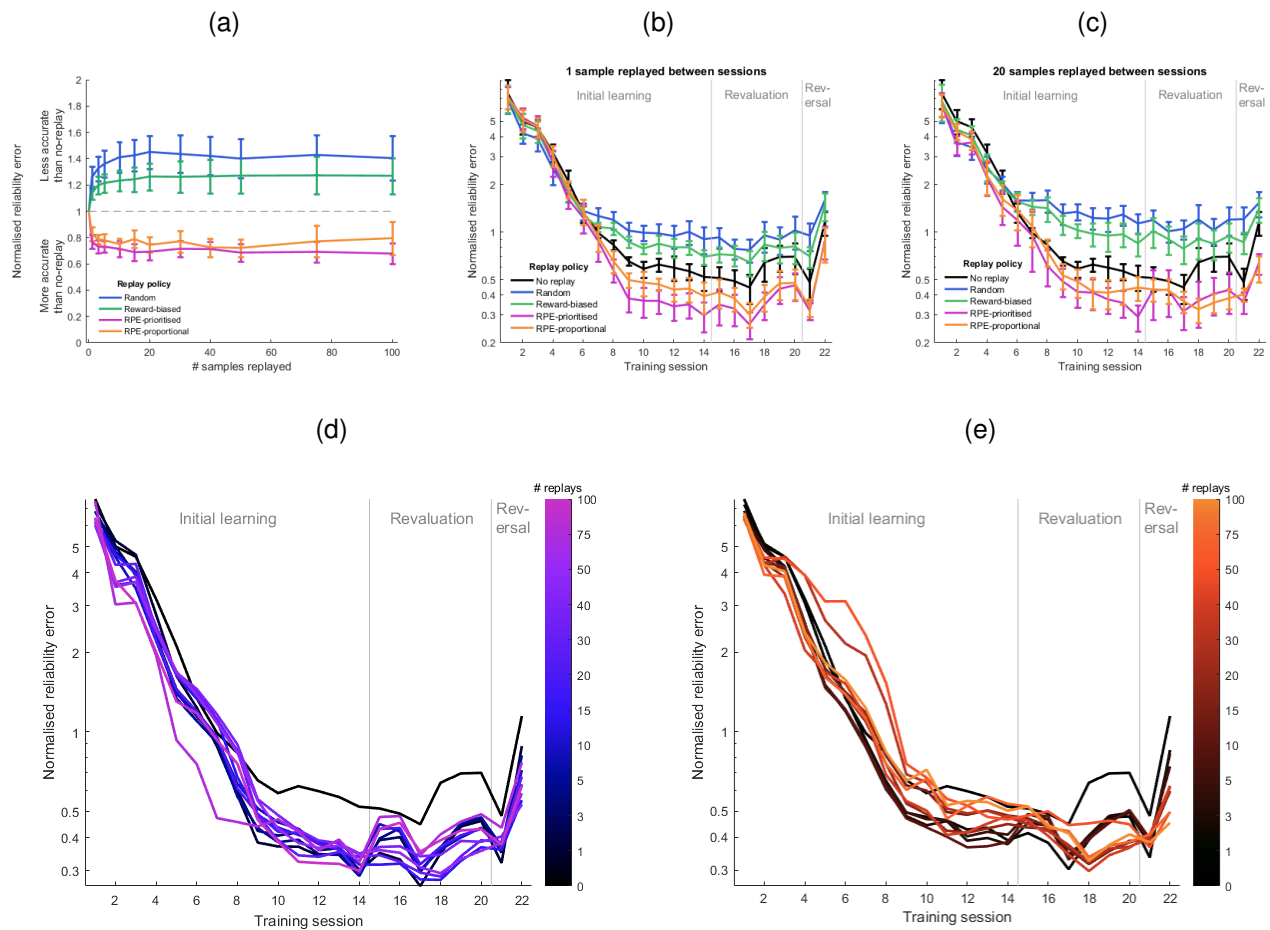


Figure 4: **A.** Normalised reliability error with varying numbers of samples replayed between sessions, averaged over all trials, according to the four replay policies shown. Reliability errors normalised to the average reliability with no replay, for each animal. Dashed line represents baseline with no replay. **B-C.** Average reliability error for each session, normalised to the average reliability error for no-replay for each animal. With 1 sample replayed between each session (B.) and 20 samples replayed between each session (C.). Error bars represent s.e.m. **D-E.** Average normalised reliability error for each session, with varying numbers of samples replayed. D. RPE-prioritised replay policy. E. RPE-proportional replay policy.

205

206 The superiority of the two RPE-biased replay policies was not uniform over the whole training period, however,
 207 and two patterns emerged. First, all replay policies showed improvements over no-replay in early sessions, but
 208 this effect disappeared in the random and reward-biased policies after roughly the seventh session. This initial
 209 superiority of all replay policies over no-replay cannot be due to replay itself because it begins in session 1,
 210 before any replay has taken place in the model; rather, it must be due to the non-replay parameters. Specifically,
 211 the optimised exploration parameter ϵ was higher in all replay policies than no-replay, so it may be the case
 212 that animals tended more towards exploration and relied on Q-values less in early training sessions. The higher
 213 ϵ value in the replay policies therefore better modelled behaviour in early sessions, whereas the differences in
 214 Q-values resulting from different replay policies impacted behaviour only later.

215 The second notable pattern is the fluctuations in the reliability errors over training sessions. In the no-replay
 216 baseline, reliability error increased in sessions 18-20 and in session 22 ($t=3.54$, $p=1.8e^{-3}$, t-test compared
 217 to reliability error in sessions 15-17 and session 21). This mirrors an increase in optimal behaviour in these

218 sessions during the revaluation stage and reversal stage respectively, suggesting that the model failed to capture
219 subtleties in the learning pattern at these points when animals were adapting their behaviour to changes in re-
220 ward probabilities. As animals re-evaluated the state-action pairs in sessions 18-20 and adjusted their behaviour
221 accordingly, replay by any policy was sufficient to overcome the increase in reliability error seen in the baseline,
222 so there was no increase at these sessions (fig. 4c; $p=0.37$ for random replay, $p=0.94$ for reward-biased replay,
223 $p=0.081$ for RPE-prioritised replay, $p=0.06$ for RPE-proportional replay with 20 samples replayed, sessions 18-
224 20 compared to sessions 15-17). This may reflect the faster learning enabled by replaying recently experienced
225 trials. However, as animals reversed their behaviour in session 22, requiring a substantial update to Q-values
226 and a dramatic change in behaviour, increased random replay or reward-biased replay did not improve reliabil-
227 ity error. With increased RPE-prioritised or RPE-proportional replay, on the other hand, increasing replay had a
228 particularly strong effect on improving reliability error in session 22 (fig. 4d-4e). This raises the possibility that
229 RPE-biased replay is especially important for behavioural flexibility of the kind seen in the reversal learning
230 stage.

231 RPE-biased replay did not improve predictions when trained on shuffled 232 data

233 Given the indication that replay might play different roles in different learning stages, it is important to control
234 for the possibility that parameter values were optimised for the general statistics of rewards and actions in the
235 task, rather than truly modelling the learning curve. Otherwise, the apparent superiority of RPE-biased replay
236 may result from anomalous irregularities in the learning patterns and not true cognitive processes. Therefore,
237 the same algorithms were trained on shuffled behavioural data in which the order of trials was randomly per-
238 muted 1,000-fold. This preserved the average frequency of state-action pairs and their associated rewards,
239 as well as the lengths of training sessions, but altered the learning curve including revaluation and reversal
240 learning.

241 Overall, the reliability errors for Q-learning with no replay were lower for shuffled data than real data, because
242 shuffled behaviour was necessarily more consistent over time and therefore more predictable. Similarly to real
243 data, reliability errors decreased sharply in early training sessions before reaching an asymptotic level (fig. 5),
244 because Q-values in early training sessions were distorted by unrepresentative rewards as a result of a small
245 sample size of trials experienced. Unlike real data, the approach to asymptotic reliability error was smooth and
246 monotonic.

247

248 Crucially, compared to the no-replay baseline, no replay policy improved reliability error. All replay policies
249 resulted in higher normalised reliability errors than no-replay ($p=6.9e^{-6}$ random, $p=6.9e^{-6}$ reward-biased,
250 $p=1.6e^{-5}$ RPE-prioritised, $p=3.4e^{-5}$). This confirms that the improvement in reliability error in the real data is
251 a result of better predictions of the learning process, and not better convergence to general statistics in the task.

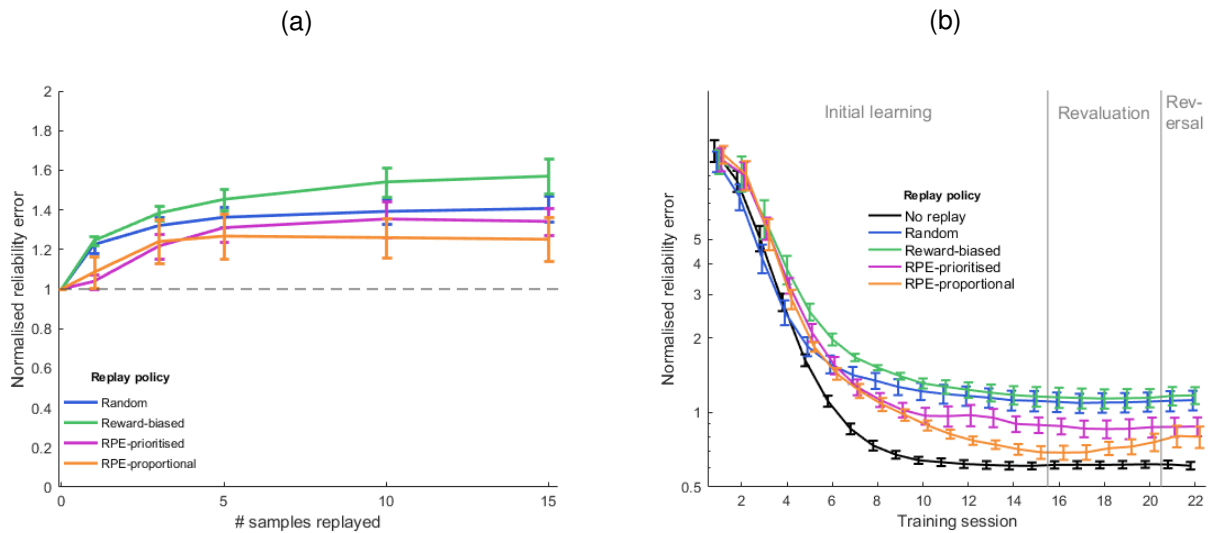


Figure 5: **A.** Normalised reliability error with varying numbers of samples replayed between sessions, trained on shuffled data in which trial data (state, action and reward) are randomly permuted. Dashed line represents baseline with no replay. **B.** Average reliability error for each session of shuffled data, normalised to the average reliability error for no-replay for each animal, with 15 samples replayed between each session. Error bars represent s.e.m.

252 **Replay-biased RPE was the best predictor for all state-action pairs**

253 We next accounted for the skew in training data towards the state-action pairs that were chosen most frequently.
254 The transition from the high-probability arm to the mid-probability arm and vice versa (as they were in the
255 initial and revaluation learning stages) were the most commonly experienced state-action pairs, representing
256 42% of trials overall, and the reliability error was weighted by the frequency of each state such that errors in
257 the more common states contributed more to the overall reliability error than errors in the less common states.
258 We therefore confirmed that Q-learning with RPE-biased replay learned to correctly predict all actions and not
259 just the more-frequently chosen actions to which the cost function was skewed.

260 Figure 6 shows the improvement in reliability errors for each replay policy over no-replay baseline, for each
261 state-action pair separately. Despite the skew in training data, the RPE-biased replay policies outperformed ran-
262 dom and reward-biased replay policies for every state-action pair, although the improvement was not identical
263 in each case. Nevertheless, the broad conclusion can be reached that RPE-biased replay policies better predicted
264 learning than either no-replay, random replay or reward-biased replay for all state-action pairs.

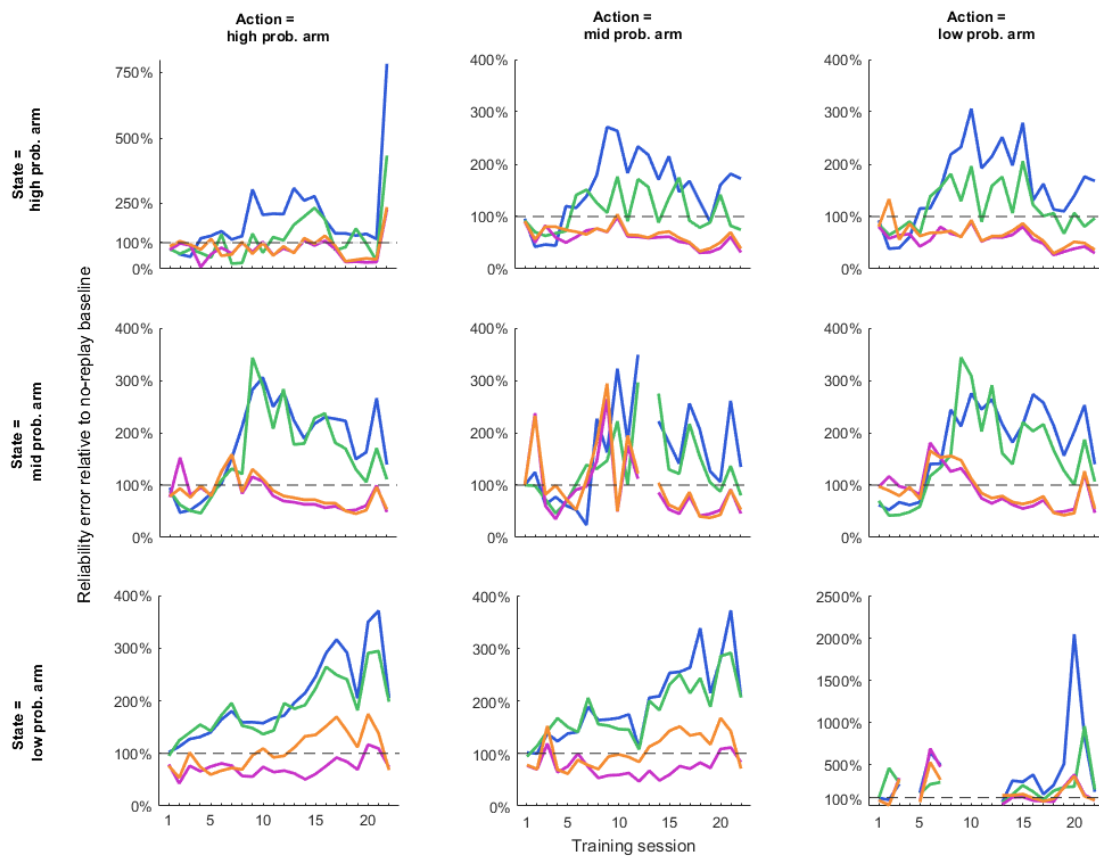


Figure 6: Change in reliability error for all trials on which a given state-action pair was expressed, with 15 samples replayed, relative to no-replay baseline. Intersection of "State = high prob. arm" and "Action = mid prob. arm" indicates a transition from high-probability arm to mid-probability arm.

265 **3 Discussion**

266 We trained rats on a reinforcement learning task designed to dissociate reward outcome (presence or absence
267 of reward) from reward prediction error (RPE; an unexpected reward or absence of reward) on each trial. We
268 trained variations of a Q-learning reinforcement learning model to predict behaviour on the task, and found that
269 Q-learning with replay prioritised by RPE was the best predictor of learning.

270 Our first main result was that Q-learning can suitably model rats' learning of the stochastic reinforcement
271 learning task, producing low reliability-errors when trained on rats' behaviour and predicting the likelihood of
272 actions on each trial. This is consistent with other studies showing that Q-learning can predict behaviour in
273 a range of tasks in rodents, monkeys and humans (Ito and Doya 2009). Given this result, we then proposed
274 that adding replay to the Q-learning model between sessions might better reflect learning and therefore better
275 predict behaviour. However, under a policy of replaying state-action pairs randomly, this produced higher
276 reliability errors overall, indicating a worse model of the cognitive processes underlying reinforcement learning.
277 Similarly, biasing replay by sampling from state-action pairs which had produced the largest recent reward also
278 increased reliability errors relative to no-replay.

279 In contrast, biasing replay by sampling from state-action pairs which had produced the largest recent RPE de-
280 creased reliability errors. From this we conclude that the cognitive processes involved in the learning of this task
281 are influenced by offline activity that takes place between sessions. Performance on memory tasks has widely
282 been found to improve following a period of sleep (Stickgold 2005; Marshall and Born 2007; Diekelmann and
283 Born 2010), associated with replay of activity which encodes recent experiences during hippocampal sharp-
284 wave ripples (Ólafsdóttir et al. 2018). We therefore propose that such offline replay underlies the RPE-biased
285 offline updating of state-action values which influenced reinforcement learning in this task.

286 The suggestion that hippocampal replay might be biased by RPEs differs from the commonly held view that
287 replay is biased by reward itself (Ambrose et al. 2016; Atherton et al. 2015; Gruber et al. 2016; Singer and
288 Frank 2009). However, the studies on which this conclusion is based generally do not use tasks which explicitly
289 dissociate reward from RPE, so these results in the literature are not inconsistent with our suggestion that RPE
290 biases replay.

291 Our conclusion that RPE-biased replay (but not random or reward-biased replay) improved model predictions
292 is strengthened by the fact that this result did not hold when training data was shuffled. When the trial order was
293 shuffled, such that there was no correlation between learning and behaviour, all replay policies produced higher
294 reliability errors in predicting the animals' behaviour. This means that the influence of RPE is a feature of the
295 learning process and not an epiphenomenon resulting from the general statistics of behaviour. Moreover, the
296 result did hold for all state-action pairs, despite the overrepresentation in training data of those most frequently
297 experienced. This gives credence to the notion that the Q-learning model with replay biased by RPE is a good
298 overall model of state-action values held by the brain.

299 Despite the prevalence of the idea that reward biases replay, our alternative theory that RPE biases replay fits
300 better with existing research on the role of dopamine. Dopaminergic projections from the ventral tegmental area
301 (VTA) to CA1 in the hippocampus have been found to modulate both replay during sleep following exposure to

302 a novel environment, and subsequent memory performance in the same environment (McNamara et al. 2014). It
303 is suggested that dopaminergic neuromodulation might tag synapses by upregulating plasticity-related proteins,
304 causing long-lasting potentiation which allows the stabilisation of the memory trace during subsequent sleep
305 and rest (Frey and Morris 1998; Redondo and Morris 2011). Phasic dopaminergic inputs to the hippocampus
306 are triggered not only in response to novelty, but also in the context of reward (Schultz et al. 1997), offering
307 a likely mechanism by which reward-related information might influence replay. Indeed, post-task replay has
308 been found in reward-related VTA cells (Gomperts et al. 2015; Valdés et al. 2015). However, such phasic
309 dopamine activations are typically elicited in response to anticipation of reward and RPEs rather than reward
310 itself (D'Ardenne et al. 2008; Dayan and Niv 2008; Montague et al. 1996; Schultz 1998; Schultz et al. 1997).
311 These phasic dopamine signals could therefore bias hippocampal replay towards activity associated with RPEs;
312 it is less clear how activity associated with reward per se might bias replay.

313 Several studies have expressly linked replay to reward, ostensibly in contrast with our results, but often RPE is
314 a confounding factor in these which cannot be discounted. In humans, high monetary reward (but not low mon-
315 etary reward) is linked to sleep-dependent improvements in associative memory (Igloi et al. 2015; Studte et al.
316 2017); in this task RPE was not estimated but would presumably be higher overall in the high-reward than low-
317 reward condition, conflating reward-dependent effects with RPE-dependent effects. In rodents, newly-rewarded
318 behaviour has been associated with replay more than behaviour which had been rewarded in previous sessions
319 (Singer and Frank 2009); here, the authors attributed this replay bias to novelty, but it is also consistent with in-
320 creased RPE when new behaviours are rewarded for the first time. Moreover, following extended reinforcement
321 of both behaviours, the replay bias for the newly-rewarded behaviour was eliminated. In a third study, results
322 were more mixed: following an increase in reward magnitude at one end of a linear track, there was more replay
323 associated with the larger-magnitude end than the unchanged-magnitude end, correlated with both reward and
324 RPE (Ambrose et al. 2016). However, following an elimination of reward at one end, there was a reduction in
325 replay following a reduction in reward despite the increase in RPE. This is more consistent with reward-biased
326 than RPE-biased replay, although the authors noted a rebound effect when the eliminated reward was reinstated:
327 greater replay was found at the reinstated-reward end than the unchanged-reward end, despite identical reward
328 magnitudes. This leaves open the possibility of bias by positive over negative RPEs. A fourth study found more
329 replay of large-reward-related activity than small-reward-related activity on a maze task (Michon et al. 2019),
330 but because reward was received on every trial analysed, any effects of reward magnitude are conflated with
331 positive reward-prediction error.

332 Conversely, the specific case for RPE-biased replay is supported by findings that neural sensitivity to RPEs in
333 humans predicts the amount of awake replay during a reinforcement learning task, and replay amount correlated
334 with subsequent performance in a task requiring behavioural flexibility (Momennejad et al. 2018).

335 In addition to human and rodent studies, findings from the literature on machine learning show some con-
336 sistency with our results. A number of machine learning studies have found that storing new information in
337 memory buffers and sampling from it at regular intervals, similar to hippocampal replay, can speed up learning
338 (Lin 1992; Mnih et al. 2013, Mnih et al. 2015), and more so when replay is biased by prediction errors (Cichosz
339 1999; Schaul et al. 2016). RPE-biased replay may therefore represent an adaptive focus whereby resources are
340 focused on areas of a cognitive model which needs updating.

341 We do not claim that this tells the whole story: RPE is almost certainly not the only factor that biases replay

342 and the phenomenon is likely to be much more multifaceted than this model suggests. First, phasic dopamine
343 signalling to hippocampus may encode other kinds of prediction errors or aspects of reward to which the VTA is
344 sensitive (Keiflin et al. 2019; Sharpe et al. 2019; Takahashi et al. 2017), and bias replay by the same mechanism.
345 Reward itself may bias replay, especially if positive RPEs influence replay more than negative RPEs; there is
346 also evidence that novelty (Hirase et al. 2001; Kudrimoti et al. 1999), the expectation of reward (Gruber et al.
347 2016), frequency of experience (Gupta et al. 2010) and strength of encoding (Schapiro et al. 2018) bias replay
348 too. Furthermore, in addition to aiding reinforcement learning, replay has been associated with other memory-
349 related functions including planning (Ólafsdóttir et al. 2017; Pfeiffer and Foster 2013), processing of emotional
350 memories (Genzel et al. 2015), creative problem-solving (Lewis et al. 2018), and generalising from episodic
351 memories to abstractions (Lewis and Durrant 2011), all of which are likely to necessitate some biasing of replay
352 distinct from RPEs. In sum, we submit that hippocampal replay is more complex than the model outlined here.

353 Our model assumes that a cache of all experience is stored from which to be sampled, which is expensive and
354 unrealistic at large scales. This may not be necessary if memory for individual trials is gradually forgotten
355 and subsumed into cortical long-term memory, for example over the course of hours over which cell assembly
356 activation decays (Giri et al. 2019).

357 Finally, this model leaves open some questions. It will be necessary to directly test this theory by recording
358 neural data from which replay can be directly observed, comparing replay of reward-associated activity with
359 that of RPE-related activity in the VTA or striatum. There is also an open question about possible diverging
360 roles of replay during behaviour compared to prolonged rest and sleep. Here we have considered replay between
361 sessions, which is likely to take place at least partly during sleep; but replay during wake has also been shown
362 to be necessary for learning (Jadhav et al. 2012).

363 In summary, we found that a Q-learning-based reinforcement learning model which assumes offline updates
364 between sessions is a better predictor of learning behaviour than one which does not assume offline updates.
365 Specifically, this is true when updates are prioritised according to experiences that have recently elicited high
366 RPEs, and not when they are prioritised according to reward or random recent experiences. This finding offers
367 a reinterpretation of how offline activity during rest and sleep might aid reinforcement learning, in terms of
368 RPE rather than reward.

369 **4 Materials and Methods**

370 **Behavioural task**

371 All procedures were performed in accordance with the United Kingdom Animals (Scientific Procedures) Act
372 1986 and European Union Directive 2010/63/EU and were reviewed by the University of Bristol Animal Wel-
373 fare and Ethical Review Board.

374 Six adult male Lister hooded rats (weighing 260-330g, Charles River Laboratories, UK) were individually
375 housed with environmental enrichment, and food-restricted to no less than 85% of their pre-restriction body
376 weight. They were trained during the light part of a 12:12 light/dark cycle to forage on a 3-armed radial
377 maze for liquid sucrose rewards in a dimly-lit room. The maze consisted of a raised central platform 25cm in
378 diameter, with three arms (60cm x 7cm) protruding from it (fig. 1a). Arms were separated from the central
379 platform by inverted-guillotine doors, which raised to block access to the arms, and fell below the maze floor to
380 allow access. Turning zones (10 x 10cm) with lick ports were positioned at the end of each arm, at which 20%
381 sucrose solution rewards were delivered. Door movements and reward delivery were operated automatically
382 according to the animal's position, tracked using a webcam mounted above the maze, using custom MATLAB
383 (The MathWorks) code. Following at least three days of habituation to the recording room and maze-operation
384 sounds, each animal performed 22 training sessions, between 5 and 7 days per week, lasting 1 hour each.

385 Trials began when a rat entered, or was placed by the experimenter on, the central platform with all doors
386 closed. Doors opened following a 5-second delay period. When the animal reached the lick port, reward was
387 probabilistically delivered or withheld, and doors to the other two arms were closed; the third door was closed
388 when the animal re-entered the central platform to begin a new trial.

389 Each arm was assigned as either "high probability", "mid probability" or "low probability", which determined
390 the protocol for reward delivery. These assignments remained fixed throughout training for each animal, but
391 were counter-balanced between animals. For the first 15 training sessions, the high-probability arm delivered
392 a reward on 6 out of 8 (75%) legitimate entries to the arm, the mid-probability arm on 4 out of 8 (50%), and
393 the low-probability arm on 2 out of 8 (25%). A legitimate entry was one in which a different arm had been
394 entered on the previous trial; entering the same arm twice in a row was incorrect and did not result in a reward
395 delivery. For sessions 16-20, the reward probabilities for the high- and low-probability arms were amplified:
396 reward was delivered on 7 out of 8 (87.5%) and 1 out of 8 (12.5%) legitimate entries respectively. For sessions
397 21-22 the reward probabilities for the high- and low-probability arms were switched, such that the (formerly)
398 high- and low- probability arms delivered reward on 1 out of 8 (12.5%) and 7 out of 8 (87.5%) of legitimate
399 entries respectively.

400 Q-learning

401 We trained several variations of a Q-learning algorithm on the behavioural data to predict choices of which
402 arm would be entered on each trial. Q-learning is a reinforcement learning algorithm developed for Markov
403 decision processes in which an agent selects actions in its environment and observes the outcome, recording at
404 each time step t its starting state s_t , selected action a_t , resulting reward r_t , and resulting state s_{t+1} . The agent
405 builds up a matrix Q of Q-value estimates for every state-action pair:

$$\begin{bmatrix} Q_{s_1,a_1} & Q_{s_1,a_2} & \cdots & Q_{s_1,a_A} \\ Q_{s_2,a_1} & Q_{s_2,a_2} & \cdots & Q_{s_2,a_A} \\ \vdots & \vdots & \ddots & \vdots \\ Q_{s_S,a_1} & Q_{s_S,a_2} & \cdots & Q_{s_S,a_A} \end{bmatrix} \quad (3)$$

406 corresponding to the future discounted expected reward, i.e. the temporal difference between the current state
407 and the reward state. These Q-value estimates are used to guide actions to maximise reward. At each time step
408 t , the Q-value for the state-action pair observed is updated by:

$$Q(s_t, a_t) \leftarrow (1 - \alpha) \cdot Q(s_t, a_t) + \alpha \cdot (r_t + \gamma \cdot \max_a Q(s_{t+1}, a)) \quad (4)$$

409 where $\alpha \in (0, 1)$ is a learning rate parameter which determines the degree to which new information overrides
410 old information, and $\gamma \in (0, 1)$ is a discount parameter which determines the importance of long-term gains.

411 In this task, entries into a chosen arm (and arrival at the goal location at the end of the arm) were modelled
412 as actions, while the arm entered on the previous trial, on which reward probabilities were contingent, were
413 modelled as states. Each trial therefore gave rise to one state-action transition out of nine possible state-action
414 pairs.

415 Q-learning with replay

416 We used four variants of Q-learning in which additional "offline" updates are performed between "online"
417 trials, based on sequences already experienced, to boost learning. This has the effect of learning from several
418 trials per actual trial of experience, and is similar to the Dyna-Q algorithm which has been shown to speed
419 up learning compared to Q-learning alone (Sutton 2014) in a manner which may underlie the function of
420 hippocampal replay (Johnson and Redish 2005). Generally, sequences are selected randomly from a memory
421 buffer of recently-acquired experiences, without bias towards any trial or type of trial. Given the observed bias
422 reported in the literature towards salient experiences, such as those rewarded or aversive, we modified Dyna-Q
423 to perform updates only between sessions and to reflect hypothesised biases in four different ways.

424 Parameter-fitting

425 Parameter-fitting for Q-learning

426 First, a Q-learning algorithm (without replay) was trained, to obtain a baseline score against which various
427 replay policies could be compared. Q-values were stored for each state-action pair on the task, and updated
428 according to each animal's experience. A state s_t was defined as the arm visited on the previous trial $t - 1$, and
429 an action a_t was defined as the arm chosen on the current trial t . Following each trial of an animal's training,
430 the Q-value $Q(s_t, a_t)$ was updated according to the reward received, $r \in \{0, 1\}$ by equation 4, and Q-values
431 were transformed into a forecast probability of choosing each arm on the subsequent trial.

432 The learning rate α , discount factor γ , and exploration factor ϵ were free parameters that were tuned to each
433 rat, using the following optimisation procedure. Here we used a reliability score (Murphy and Murphy 1973),
434 generated based on the forecast probabilities of all trials, to quantify the consistency of the forecast probabilities
435 with the animals' behaviour. The mean observed frequency was calculated for each state-action pair, i.e. the
436 proportion of trials on which a given action was chosen in a given state, and the reliability score R_t for a given
437 trial t was calculated according to:

$$R_t = n_{s_t} \cdot \sum_{a=1}^{n_a} (p_a - o_{s_t,a})^2 \quad (5)$$

438 where s_t is the animal's state on trial t , n_{s_t} is the number of trials on which the animal was in state s_t , n_a is the
439 number of possible actions (3) p_a is the forecast probability for entering arm a , and $o_{s,a}$ is the mean observed
440 frequency of state-action pair s, a .

441 Parameter optimisation was performed using the reliability error as the cost function. Because the parameter
442 state-space was vulnerable to local minima, and also because it was highly stochastic under replay policies (see
443 description below), a two-step approach was taken to optimise parameters. In the first step, simulated annealing
444 was run 32 times for a maximum of 1000 iterations (or until the reliability error could not be improved by more
445 than $1e^{-6}$), using the MATLAB function `simulannealbnd`. Reliability error was averaged over 1,000 runs
446 when computing the cost function, to minimise stochasticity. This function performs a probabilistic variation of
447 gradient descent by taking increasingly smaller steps in random directions, to approximate a global minimum
448 without becoming stuck in local minima. The resulting 32 rough estimates of the optimal parameter values
449 were used as the initial values for the second step: a simple quasi-Newton method using gradient descent,
450 implemented by the MATLAB optimisation function `fmincon`, for a further maximum 1000 iterations. Of the
451 32 final sets of parameter values, the one which produced the smallest reliability error was used for analysis.
452 All analyses were performed on the average reliability error over 1,000 runs using the given parameter values.

453 **Parameter-fitting for Q-learning with replay**

454 Against the baseline of no-replay, the same optimisation procedure was performed with increasing amounts of
455 replay according to four replay policies. Following each session, a specified number of samples were chosen
456 from all the trials experienced so far. How the samples were selected depended on the replay policy (detailed
457 below); a probability $P(s, a)$ was assigned to each state-action pair to determine which pair to sample from.
458 From the chosen state-action pair, a sample trial was chosen according to the probability $P(i)$ in which a
459 recency parameter ensured that more recent trials were exponentially more likely to be chosen. Q-values were
460 then updated according to the state, action and reward of the sampled trial, in the same manner as "online"
461 Q-value updates described in equation 4.

462 Each replay policy required the same three parameters to be optimised as in Q-learning without replay, plus
463 additional parameters for recency and/or RPE-weighting. Table 2 shows the number of free parameters for each
464 replay policy.

Replay policy	Number of parameters
No replay	3
Random replay	4
Reward-biased replay	4
RPE-prioritised replay	5
RPE-proportional replay	5

Table 2

465 These were optimised according to the same procedure as for Q-learning with no replay, described above,
466 for $n = \{1, 3, 5, 10, 15, 20, 30, 40, 50, 75, 100\}$ replay events between each session, resulting in 11 sets of
467 parameter values for each replay policy and each animal. Comparing this to plausible quantities of replay
468 events in animals is not trivial, but studies in which discrete replay events are enumerated report 100-200 bursts
469 of hippocampal activity that can be statistically related to prior experience, over the first one or two hours after
470 experience (Ólafsdóttir et al. 2016; Michon et al. 2019). Separately, reactivation of cell pairs has been found to
471 decay to baseline well within that time period following exposure to familiar environments (Giri et al. 2019),
472 so the first one to two hours is likely to be when most replay of recent experience in a familiar environment
473 occurs.

474 **Random replay**

475 Random replay, biased by nothing but the recency of an action, was included as a control. For each replay
476 event, a state-action pair was chosen at random out of all state-action pairs experienced so far:

$$P(s, a) = \frac{1}{n_{sa}} \quad (6)$$

477 where n_{sa} is the number of state-action pairs experienced (up to 9). The subset of trials experienced, $i \in (1, I)$,
478 which represented this state-action pair were ordered chronologically, and the probability $P(i)$ of a trial i being
479 replayed was determined according to a recency parameter φ :

$$P(i) = \frac{i^\varphi}{\sum_{i=1}^I p_i} \quad (7)$$

480 **Reward-biased replay**

481 Reward-biased replay represents the predominant interpretation of how reward influences replay (Atherton et al.
482 2015, Carr et al. 2011). For each replay event, a state-action pair s, a was chosen probabilistically in proportion
483 to its Q-value:

$$P(s, a) = \frac{Q(s, a)}{\sum_{s=1}^{n_s} \sum_{a=1}^{n_a} Q(s, a)} \quad (8)$$

484 The subset of trials experienced which represented the chosen state-action pair were ordered chronologically,
485 and determined according to equation 7.

486 **RPE-prioritised replay**

487 RPE-prioritised replay represents the policy of replaying trials associated with the most surprising outcomes,
488 i.e. where the difference between expectation (Q-values) and experience (reward) was greatest. For each trial
489 t , RPE was calculated as the difference between actual reward and expected reward:

$$\text{rpe}_t = r + \gamma \cdot Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t) \quad (9)$$

490 For every trial $i \in (1, I)$ which was an example of a given state-action pair, its absolute value was weighted,
491 determined by a parameter φ raised to the power of its recency i :

$$\text{wrpe}_i = |\text{rpe}_i| \cdot \varphi^i \quad (10)$$

492 The weighted RPEs, wrpe , were then averaged to produce an overall weighted-average RPE, $\text{RPE}_{s,a}$, for each
493 state-action pair s, a , which was more heavily influenced by recent trials:

$$\text{RPE}_{s,a} = \frac{\sum_{i=1}^I \text{wrpe}_i}{I} \quad (11)$$

494 The state-action pair with the highest RPE was selected, and the subset of trials experienced which represented
495 the chosen pair were ordered chronologically, and determined according to equation 7. Once replayed, the rpe_t
496 for the trial sampled was updated to reflect the RPE resulting from the replay event.

497 **RPE-proportional replay**

498 RPE-proportional replay is a variant of RPE-prioritised replay, in which state-action pairs are chosen in propor-
499 tion to their weighted-average-RPE instead of choosing the pair with the highest weighted-average-RPE. The
500 RPE was calculated according to eq. 11 and a state-action pair to be sampled from was chosen probabilistically
501 according to:

$$p_{s,a} = \frac{\text{RPE}_{s,a}}{\sum \text{RPE}_{s,a}} \quad (12)$$

502 The subset of trials experienced which represented the chosen state-action pair were ordered chronologically,
503 and determined according to equation 7. Once replayed, the rpe_t for the trial sampled was updated to reflect
504 the RPE resulting from the replay event.

505 **Shuffling procedure**

506 As an additional control, the parameters were also optimised for shuffled data, in which trial order was randomly
507 permuted 1,000-fold. This preserved the large-scale information in the training data, such as the mean observed
508 frequency and average rewards of state-action pairs and the number of trials in each session between replays,
509 but disrupted the specific structure of how this information was acquired over time.

510 **Code Availability**

511 All data and code used in this study are available at <https://github.com/eroscow/QlearningReplay>.

512 References

- 513 Ambrose, R. E., Pfeiffer, B. E., Correspondence, D. J. F. & Foster, D. J. (2016). Reverse Replay of
514 Hippocampal Place Cells Is Uniquely Modulated by Changing Reward. *Neuron*, *91*. doi:10.1016/
515 j.neuron.2016.07.047
- 516 Andrychowicz, M., Wolski, F., Ray, A., Schneider, J., Fong, R., Welinder, P., . . . Zaremba, W. (2017).
517 Hindsight Experience Replay. Retrieved from [http://papers.nips.cc/paper/7090-hindsight-experience-](http://papers.nips.cc/paper/7090-hindsight-experience-replay)
518 [replay](http://papers.nips.cc/paper/7090-hindsight-experience-replay)
- 519 Antony, J. W., Gobel, E. W., O'Hare, J. K., Reber, P. J. & Paller, K. A. (2012). Cued memory reactiva-
520 tion during sleep influences skill learning. *Nat. Neurosci.* *15*(8), 1114–6. doi:10.1038/nn.3152
- 521 Atherton, L. A., Dupret, D. & Mellor, J. R. (2015). Memory trace replay: the shaping of memory
522 consolidation by neuromodulation. *Trends Neurosci.* *38*(9), 560–70. doi:10.1016/j.tins.2015.07.
523 004
- 524 Bendor, D. & Wilson, M. A. (2012). Biasing the content of hippocampal replay during sleep. *Nat.*
525 *Neurosci.* *15*(10), 1439–44. doi:10.1038/nn.3203
- 526 Cairney, S. A., Durrant, S. J., Jackson, R. & Lewis, P. A. (2014). Sleep spindles provide indirect
527 support to the consolidation of emotional encoding contexts. *Neuropsychologia*, *63*, 285–292.
528 doi:10.1016/j.neuropsychologia.2014.09.016
- 529 Calabresi, P., Picconi, B., Tozzi, A. & Di Filippo, M. (2007). Dopamine-mediated regulation of cor-
530 ticostriatal synaptic plasticity. *Trends Neurosci.* *30*(5), 211–219. doi:10.1016/J.TINS.2007.03.001
- 531 Carr, M. F., Jadhav, S. P. & Frank, L. M. (2011). Hippocampal replay in the awake state: a potential
532 substrate for memory consolidation and retrieval. *Nat. Neurosci.* *14*(2), 147–53. doi:10.1038/nn.
533 2732
- 534 Cheng, S. & Frank, L. M. (2008). New experiences enhance coordinated neural activity in the hippo-
535 campus. *Neuron*, *57*(2), 303–13. doi:10.1016/j.neuron.2007.11.035
- 536 Cichosz, P. (1999). An analysis of experience replay in temporal difference learning. *Cybern. Syst.*
537 *30*(5), 341–363. doi:10.1080/019697299125127
- 538 D'Ardenne, K., McClure, S. M., Nystrom, L. E. & Cohen, J. D. (2008). BOLD responses reflecting
539 dopaminergic signals in the human ventral tegmental area. *Science*, *319*(5867), 1264–7. doi:10.
540 1126/science.1150605
- 541 Daw, N. D., Niv, Y. & Dayan, P. (2005). Uncertainty-based competition between prefrontal and dor-
542 solateral striatal systems for behavioral control. *Nat. Neurosci.* *8*(12), 1704–1711. doi:10.1038/
543 nn1560
- 544 Day, J. J., Roitman, M. F., Wightman, R. M. & Carelli, R. M. (2014). Associative learning mediates
545 dynamic shifts in DORamine signaling in the nucleus accumbens Associative learning mediates
546 dynamic shifts in dopamine signaling in the nucleus accumbens. doi:10.1038/nn1923
- 547 Dayan, P. & Niv, Y. (2008). Reinforcement learning: The Good, The Bad and The Ugly. *Curr. Opin.*
548 *Neurobiol.* *18*(2), 185–196. doi:10.1016/J.CONB.2008.08.003
- 549 Diekelmann, S. & Born, J. (2010). The memory function of sleep. *Nat. Rev. Neurosci.* *11*(2), 114–126.
550 doi:10.1038/nrn2762
- 551 Dupret, D., O'Neill, J. & Pleydell-Bouverie, B. (2010). The reorganization and reactivation of hippo-
552 campal maps predict spatial memory performance. *Nature*. Retrieved from [http://www.nature.](http://www.nature.com/neuro/journal/v13/n8/abs/nn.2599.html)
553 [com/neuro/journal/v13/n8/abs/nn.2599.html](http://www.nature.com/neuro/journal/v13/n8/abs/nn.2599.html)

- 554 Ego-Stengel, V. & Wilson, M. A. (2010). Disruption of ripple-associated hippocampal activity during
555 rest impairs spatial learning in the rat. *Hippocampus*, *20*(1), 1–10. doi:10.1002/hipo.20707
- 556 Foster, D. J. & Wilson, M. A. (2006). Reverse replay of behavioural sequences in hippocampal place
557 cells during the awake state. *Nature*, *440*(7084), 680–683. doi:10.1038/nature04587. arXiv: 440:
558 680–683
- 559 Frey, U. & Morris, R. G. (1998). Synaptic tagging: implications for late maintenance of hippocampal
560 long-term potentiation. *Trends Neurosci.* *21*(5), 181–188. doi:10.1016/S0166-2236(97)01189-2
- 561 Genzel, L., Spormaker, V., Konrad, B. & Dresler, M. (2015). The role of rapid eye movement sleep
562 for amygdala-related memory processing. *Neurobiol. Learn. Mem.* *122*, 110–121. doi:10.1016/
563 J.NLM.2015.01.008
- 564 Girardeau, G., Benchenane, K., Wiener, S. I., Buzsáki, G. & Zugaro, M. B. (2009). Selective sup-
565 pression of hippocampal ripples impairs spatial memory. *Nat. Neurosci.* *12*(10), 1222–1223.
566 doi:10.1038/nn.2384
- 567 Girardeau, G., Inema, I. & Buzsáki, G. (2017). Reactivations of emotional memory in the hippocam-
568 pus–amygdala system during sleep. *Nat. Neurosci.* *20*(11), 1634–1642. doi:10.1038/nn.4637
- 569 Giri, B., Miyawaki, H., Mizuseki, K., Cheng, S. & Diba, K. (2019). Hippocampal Reactivation Ex-
570 tends for Several Hours Following Novel Experience. *J. Neurosci.* *39*(5), 866–875. doi:10.1523/
571 JNEUROSCI.1950-18.2018
- 572 Gomperts, S. N., Kloosterman, F., Wilson, M. A., Cardinal, R., Parkinson, J., Hall, J., . . . Sejnowski, T.
573 (2015). VTA neurons coordinate with the hippocampal reactivation of spatial experience. *Elife*,
574 *4*, 321–352. doi:10.7554/eLife.05360
- 575 Gruber, M. J., Ritchey, M., Wang, S.-F., Doss, M. K. & Ranganath, C. (2016). Post-learning Hippo-
576 campal Dynamics Promote Preferential Retention of Rewarding Events. *Neuron*, *89*(5), 1110–
577 1120. doi:10.1016/j.neuron.2016.01.017
- 578 Gupta, A. S., van der Meer, M. A., Touretzky, D. S. & Redish, A. D. (2010). Hippocampal Replay Is Not
579 a Simple Function of Experience. *Neuron*, *65*(5), 695–705. doi:10.1016/J.NEURON.2010.01.034
- 580 Hirase, H., Leinekugel, X., Czurkó, A., Csicsvari, J. & Buzsáki, G. (2001). Firing rates of hippocam-
581 pal neurons are preserved during subsequent sleep episodes and modified by novel awake
582 experience. *Proc. Natl. Acad. Sci. U. S. A.* *98*(16), 9386–90. doi:10.1073/pnas.161274398
- 583 Igloi, K., Gaggioni, G., Sterpenich, V. & Schwartz, S. (2015). A nap to recap or how reward regulates
584 hippocampal-prefrontal memory networks during daytime sleep in humans. doi:10.7554/eLife.
585 07903.001
- 586 Ito, M. & Doya, K. (2009). Validation of Decision-Making Models and Analysis of Decision Variables
587 in the Rat Basal Ganglia. *J. Neurosci.* *29*(31), 9861–9874. doi:10.1523/jneurosci.6157-08.2009
- 588 Jadhav, S. P., Kemere, C., German, P. W. & Frank, L. M. (2012). Awake Hippocampal Sharp-Wave
589 Ripples Support Spatial Memory. *Science (80-.)*. *336*(6087), 1454–1458. doi:10.1126/SCIENCE.
590 1217230
- 591 Johnson, A. & Redish, A. D. (2005). Hippocampal replay contributes to within session learning in a
592 temporal difference reinforcement learning model. *Neural Networks*, *18*(9), 1163–1171. doi:10.
593 1016/J.NEUNET.2005.08.009
- 594 Karimpanal, T. G. & Bouffanais, R. (2017). Experience Replay Using Transition Sequences. *Front.*
595 *Neurobot.* *12*, 32. doi:10.3389/fnbot.2018.00032. arXiv: 1705.10834

- 596 Keiflin, R., Pribut, H. J., Shah, N. B. & Janak, P. H. (2019). Ventral Tegmental Dopamine Neurons
597 Participate in Reward Identity Predictions. *Curr. Biol.* 29(1), 93–103.e3. doi:10.1016/J.CUB.2018.
598 11.050
- 599 Kim, H., Lee, D. & Jung, M. W. (2013). Signals for Previous Goal Choice Persist in the Dorsomedial,
600 but Not Dorsolateral Striatum of Rats. *J. Neurosci.* 33(1), 52–63. doi:10.1523/jneurosci.2422-
601 12.2013
- 602 Kudrimoti, H. S., Barnes, C. A. & McNaughton, B. L. (1999). Reactivation of hippocampal cell assem-
603 blies: effects of behavioral state, experience, and EEG dynamics. *J. Neurosci.* 19(10), 4090–
604 101. doi:10.1523/JNEUROSCI.19-10-04090.1999
- 605 Lansink, C. S., Goltstein, P. M., Lankelma, J. V., McNaughton, B. L. & Pennartz, C. M. A. (2009).
606 Hippocampus Leads Ventral Striatum in Replay of Place-Reward Information. *PLoS Biol.* 7(8),
607 e1000173. doi:10.1371/journal.pbio.1000173
- 608 Lewis, P. A. & Durrant, S. J. (2011). Overlapping memory replay during sleep builds cognitive schemata.
609 *Trends Cogn. Sci.* 15(8), 343–351. doi:10.1016/j.tics.2011.06.004
- 610 Lewis, P. A., Knoblich, G. & Poe, G. (2018). How Memory Replay in Sleep Boosts Creative Problem-
611 Solving. *Trends Cogn. Sci.* 22(6), 491–503. doi:10.1016/J.TICS.2018.03.009
- 612 Lin, L.-J. (1992). Self-improving reactive agents based on reinforcement learning, planning and teach-
613 ing. *Mach. Learn.* 8(3-4), 293–321. doi:10.1007/BF00992699
- 614 Marshall, L. & Born, J. (2007). The contribution of sleep to hippocampus-dependent memory consol-
615 idation. *Trends Cogn. Sci.* 11(10), 442–450. doi:10.1016/J.TICS.2007.09.001
- 616 McClure, S. M., Berns, G. S. & Montague, P. (2003). Temporal Prediction Errors in a Passive Learning
617 Task Activate Human Striatum. *Neuron*, 38(2), 339–346. doi:10.1016/S0896-6273(03)00154-5
- 618 McNamara, C. G., Tejero-Cantero, Á., Trouche, S., Campo-Urriza, N. & Dupret, D. (2014). Dopaminer-
619 gic neurons promote hippocampal reactivation and spatial memory persistence. *Nat. Neurosci.*
620 17(12), 1658–1660. doi:10.1038/nn.3843
- 621 Michon, F., Sun, J.-J., Kim, C. Y., Ciliberti, D. & Kloosterman, F. (2019). Post-learning Hippocampal
622 Replay Selectively Reinforces Spatial Memory for Highly Rewarded Locations. *Curr. Biol.* 29(9),
623 1436–1444.e5. doi:10.1016/J.CUB.2019.03.048
- 624 Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D. & Riedmiller, M. (2013).
625 Playing Atari with Deep Reinforcement Learning. arXiv: 1312.5602. Retrieved from [http://arxiv.
626 org/abs/1312.5602](http://arxiv.org/abs/1312.5602)
- 627 Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., ... Hassabis, D.
628 (2015). Human-level control through deep reinforcement learning. doi:10.1038/nature14236
- 629 Momennejad, I., Otto, R., Daw, N. D. & Norman, K. A. (2018). Offline replay supports planning in
630 human reinforcement learning. doi:10.7554/eLife.32548.001
- 631 Montague, P. R., Dayan, P. & Sejnowski, T. J. (1996). A framework for mesencephalic dopamine
632 systems based on predictive Hebbian learning. *J. Neurosci.* 16(5), 1936–47. doi:10.1523/
633 JNEUROSCI.16-05-01936.1996
- 634 Morris, G., Schmidt, R. & Bergman, H. (2010). Striatal action-learning based on dopamine concen-
635 tration. *Exp. Brain Res.* 200(3-4), 307–317. doi:10.1007/s00221-009-2060-6
- 636 Murphy, A. H. & Murphy, A. H. (1973). A New Vector Partition of the Probability Score. *J. Appl.*
637 *Meteorol.* 12(4), 595–600. doi:10.1175/1520-0450(1973)012<0595:ANVPOT>2.0.CO;2

- 638 O'Doherty, J. P., Dayan, P., Friston, K., Critchley, H. & Dolan, R. J. (2003). Temporal Difference Models
639 and Reward-Related Learning in the Human Brain. *Neuron*, 38(2), 329–337. doi:10.1016/S0896-
640 6273(03)00169-7
- 641 Ólafsdóttir, H. F., Bush, D. & Barry, C. (2018). The Role of Hippocampal Replay in Memory and
642 Planning. *Curr. Biol.* 28(1), R37–R50. doi:10.1016/j.cub.2017.10.073
- 643 Ólafsdóttir, H. F., Carpenter, F. & Barry, C. (2016). Coordinated grid and place cell replay during rest.
644 *Nat. Neurosci.* 19(6), 792–794. doi:10.1038/nn.4291
- 645 Ólafsdóttir, H. F., Carpenter, F. & Barry, C. (2017). Task Demands Predict a Dynamic Switch in the
646 Content of Awake Hippocampal Replay. *Neuron*, 96(4), 925–935.e6. doi:10.1016/J.NEURON.
647 2017.09.035
- 648 Pagnoni, G., Zink, C. F., Montague, P. R. & Berns, G. S. (2002). Activity in human ventral striatum
649 locked to errors of reward prediction. *Nat. Neurosci.* 5(2), 97–98. doi:10.1038/nn802
- 650 Pfeiffer, B. E. & Foster, D. J. (2013). Hippocampal place-cell sequences depict future paths to re-
651 membered goals. *Nature*, 497(7447), 74–9. doi:10.1038/nature12112
- 652 Rasch, B., Büchel, C., Gais, S. & Born, J. (2007). Odor cues during slow-wave sleep prompt declar-
653 ative memory consolidation. *Science (80-)*. Retrieved from [http://science.sciencemag.org/content/
654 315/5817/1426.short](http://science.sciencemag.org/content/315/5817/1426.short)
- 655 Redondo, R. L. & Morris, R. G. M. (2011). Making memories last: the synaptic tagging and capture
656 hypothesis. *Nat. Rev. Neurosci.* 12(1), 17–30. doi:10.1038/nrn2963
- 657 Roesch, M. R., Calu, D. J. & Schoenbaum, G. (2007). Dopamine neurons encode the better option in
658 rats deciding between differently delayed or sized rewards. *Nat. Neurosci.* 10(12), 1615–1624.
659 doi:10.1038/nn2013
- 660 Rudoy, J., Voss, J., Westerberg, C. & Paller, K. (2009). Strengthening individual memories by react-
661 ivating them during sleep. *Science (80-)*. Retrieved from [http://science.sciencemag.org/content/
662 326/5956/1079.short](http://science.sciencemag.org/content/326/5956/1079.short)
- 663 Schapiro, A. C., McDevitt, E. A., Rogers, T. T., Mednick, S. C. & Norman, K. A. (2018). Human
664 hippocampal replay during rest prioritizes weakly learned information and predicts memory
665 performance. *Nat. Commun.* 9(1), 3920. doi:10.1038/s41467-018-06213-1
- 666 Schaul, T., Quan, J., Antonoglou, I., Silver, D. & Deepmind, G. (2016). *Prioritized Experience Replay*.
667 arXiv: 1511.05952v4. Retrieved from <https://arxiv.org/pdf/1511.05952.pdf>
- 668 Schultz, W., Dayan, P. & Montague, P. R. (1997). A neural substrate of prediction and reward. *Science*,
669 275(5306), 1593–9. doi:10.1126/SCIENCE.275.5306.1593
- 670 Schultz, W. (1998). Predictive Reward Signal of Dopamine Neurons. *J. Neurophysiol.* 80(1), 1–27.
671 doi:10.1152/jn.1998.80.1.1
- 672 Schultz, W. (2016). Dopamine reward prediction error coding. *Dialogues Clin. Neurosci.* 18(1), 23–
673 32. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/27069377>
674 [http://www.pubmedcentral.
675 nih.gov/articlerender.fcgi?artid=PMC4826767](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4826767)
- 675 Sharpe, M. J., Batchelor, H. M., Mueller, L. E., Chang, C. Y., Maes, E. J., Niv, Y. & Schoenbaum, G.
676 (2019). Dopamine transients delivered in learning contexts do not act as model-free prediction
677 errors. *bioRxiv*, 574541. doi:10.1101/574541
- 678 Singer, A. C. & Frank, L. M. (2009). Rewarded Outcomes Enhance Reactivation of Experience in the
679 Hippocampus. *Neuron*, 64(6), 910–921. doi:10.1016/j.neuron.2009.11.016
- 680 Stickgold, R. (2005). Sleep-dependent memory consolidation. *Nature*, 437(7063), 1272–1278. doi:10.
681 1038/nature04286

- 682 Studte, S., Bridger, E. & Mecklinger, A. (2017). Sleep spindles during a nap correlate with post sleep
683 memory performance for highly rewarded word-pairs. *Brain Lang.* 167, 28–35. doi:10.1016/J.
684 BANDL.2016.03.003
- 685 Sutton, R. S. (2014). Integrated Architectures for Learning, Planning, and Reacting Based on Approx-
686 imating Dynamic Programming. *Mach. Learn. Proc. 1990*, 216–224. doi:10.1016/b978-1-55860-
687 141-3.50030-4
- 688 Takahashi, Y. K., Batchelor, H. M., Liu, B., Khanna, A., Morales, M. & Schoenbaum, G. (2017).
689 Dopamine Neurons Respond to Errors in the Prediction of Sensory Features of Expected Re-
690 wards Article Dopamine Neurons Respond to Errors in the Prediction of Sensory Features of
691 Expected Rewards. *Neuron*, 95. doi:10.1016/j.neuron.2017.08.025
- 692 Valdés, J. L., McNaughton, B. L. & Fellous, J.-M. (2015). Offline reactivation of experience-dependent
693 neuronal firing patterns in the rat ventral tegmental area. *J. Neurophysiol.* 114(2), 1183–95.
694 doi:10.1152/jn.00758.2014
- 695 Watkins, C. J. (1989). Learning form delayed rewards. *Ph. D. thesis, King's Coll. Univ. Cambridge.*
696 Retrieved from <https://ci.nii.ac.jp/naid/10007782517/>
- 697 Wimmer, G. E., Li, J. K., Gorgolewski, K. J. & Poldrack, R. A. (2018). Reward Learning over Weeks
698 Versus Minutes Increases the Neural Representation of Value in the Human Brain. *J. Neurosci.*
699 38(35), 7649–7666. doi:10.1523/jneurosci.0075-18.2018
- 700 Wu, C.-T., Haggerty, D., Kemere, C. & Ji, D. (2017). Hippocampal awake replay in fear memory
701 retrieval. *Nat. Neurosci.* 20(4), 571–580. doi:10.1038/nn.4507
- 702 Yu, J. Y., Kay, K., Liu, D. F., Grossrubatscher, I., Loback, A., Sosa, M., . . . Frank, L. M. (2017). Dis-
703 tinct hippocampal-cortical memory representations for experiences associated with movement
704 versus immobility. *Elife*, 6. doi:10.7554/eLife.27621

705

706 **Acknowledgements:** We are grateful to Rui Ponte Costa and Mark Humphries for useful discussions and com-
707 ments, and to Aleksander Domanski and Andrew New for assistance with experimental set-up. This research
708 was funded by a Wellcome Trust PhD scholarship.

709 **Author Contributions:** E.L.R., M.W.J. and N.F.L. conceived and designed the study. E.L.R. carried out the
710 experiments and analysed the data. E.L.R. performed the computational modelling under the guidance of N.F.L.
711 E.L.R. and N.F.L. prepared the paper with critical revision from M.W.J.

712 **Competing interests:** The authors declare no competing interests.