

# 1 Differential replay of reward and punishment paths predicts 2 approach and avoidance

3 Jessica McFadyen<sup>\*1,2</sup>, Yunzhe Liu<sup>3,4</sup>, & Raymond J Dolan<sup>1,2</sup>

4 \* corresponding author: [drjessicajeon@gmail.com](mailto:drjessicajeon@gmail.com)

5 <sup>1</sup> The UCL Max Planck Centre for Computational Psychiatry and Ageing Research, University College London, London, UK

6 <sup>2</sup> Wellcome Centre for Human Neuroimaging, University College London, London, UK

7 <sup>3</sup> State Key Laboratory of Cognitive Neuroscience and Learning, IDG/McGovern Institute for Brain Research, Beijing Normal  
8 University, Beijing, China.

9 <sup>4</sup> Chinese Institute for Brain Research, Beijing, China.

## 10 Abstract

11 Neural replay is implicated in planning, where states relevant to a task goal are rapidly reactivated in sequence. It  
12 remains unclear whether, during planning, replay relates to an actual prospective choice. Here, using  
13 magnetoencephalography (MEG), we studied replay in human participants while they planned to either approach  
14 or avoid an uncertain environment containing paths leading to reward or punishment. We find evidence for  
15 forward sequential replay during planning, with rapid state-to-state transitions from 20 to 90 ms. Replay of  
16 rewarding paths was boosted, relative to aversive paths, prior to a decision to avoid and attenuated prior to a  
17 decision to approach. A trial-by-trial bias towards replaying prospective punishing paths predicted irrational  
18 decisions to approach riskier environments, an effect more pronounced in participants with higher trait anxiety.  
19 The findings indicate a coupling of replay with planned behaviour, where replay prioritises an online representation  
20 of a worst-case scenario for approaching or avoiding.

## 21 Introduction

22 When formulating a plan, we often face uncertainty as to whether a choice will lead to a good or bad outcome.  
23 For example, when we deliberate whether to go to a party or stay home, we might simulate potential sequences  
24 of events that are positive (e.g., arriving and seeing friends, meeting new people, coming home feeling happy) or  
25 negative (e.g., arriving and not knowing anybody, saying something embarrassing in front of new people, leaving  
26 early, and regretting the whole experience). Situations such as these can engender approach-avoidance conflict,  
27 wherein decision-making is rendered difficult by a need to weigh the benefits of a risky choice against a more  
28 certain, but less rewarding, choice to avoid.

29 Neural replay, originally characterised in the context of a rapid sequential reactivation of hippocampal place cells  
30 that map specific locations of recently experienced paths<sup>1-7</sup>, is linked to a number of functions in both humans  
31 and animals, including memory consolidation of spatial<sup>4,8-11</sup> and temporal order relationships<sup>7,12,13</sup>, inference<sup>14-16</sup>,  
32 and credit assignment<sup>17</sup>. There is also evidence indicating that neural replay may relate to a simulation of potential  
33 outcomes during active planning<sup>18-22</sup>.

34 A role for prospective replay in planning is supported by observations that when rodents pause during spatial  
35 navigation, the order of replayed place cell firing matches paths leading to the learned location of a reward<sup>23,24</sup>,  
36 and is enhanced for paths leading to greater rewards<sup>25</sup>. Furthermore, the more a rewarding path is prospectively  
37 replayed, the more likely it is that the animal will pursue that path<sup>21,23,24,26</sup>. A disruption of replay events at decision  
38 points, such as by application of electric pulses to the hippocampus, leads to the expression of more vicarious  
39 trial and error behaviour<sup>25,27</sup> and a greater likelihood that an animal will take an incorrect path<sup>10,28</sup>. Remarkably,  
40 replay events also provide a mapping of potential trajectories to rewards that have never been experienced,  
41 evident in both online<sup>29</sup> and offline<sup>30,31</sup> sequential reactivation.

42 In contrast to reward, the question of how prospective aversive events modulate replay is under-investigated.  
43 Animal studies show that removal of a reward leads to a marked reduction in replay<sup>32</sup>. Paths leading to danger,  
44 however, are also more strongly replayed, and this is anticorrelated with an animal's chosen trajectory such that  
45 they tend to avoid the dangerous path<sup>33</sup>. Such findings have led to a proposal that hippocampal replay prioritises  
46 paths that are most immediately relevant for on-going behaviour<sup>34</sup>. Recent evidence, however, suggests the goal  
47 of a current plan might not, in fact, directly relate to which path is most strongly replayed. Instead, the selection  
48 of paths for replay appears to relate to mnemonic functions that support future planning, evident in replay being  
49 enhanced for paths leading to previously-rewarded locations that have not been visited recently<sup>35,36</sup>, as well as for  
50 paths associated with sub-optimal decision-making<sup>22</sup>. Within this formulation, replay is proposed to support  
51 planning by consolidating memories of sequences that are susceptible to being forgotten, rather than reflecting  
52 a simulation of states leading to outcomes that directly relate to a current motivational goal<sup>37</sup>.

53 A feature of many previous studies of replay has been the use of environments that contain either reward or  
54 punishment. Little is known about how replay is impacted by a prospective environment where paths can lead to  
55 either reward or punishment, especially where these environments give rise to an approach-avoidance conflict<sup>38,39</sup>.  
56 Notably, an inability to make optimal decisions under approach-avoidance conflict is a characteristic of clinical  
57 anxiety disorders, where the potential for experiencing a negative event leads to avoidance regardless of the  
58 likelihood of potential reward<sup>40-42</sup>. On the other hand, a tendency to approach, even when this might have negative  
59 consequences, is considered a risk factor for developing substance abuse disorders<sup>43</sup>. During approach-  
60 avoidance conflict, the magnitude and likelihood of threat is proposed to be monitored by anterior and ventral  
61 hippocampus interactions that arbitrate decisions to approach or avoid, in both humans<sup>39,44</sup> and rodents<sup>45,46</sup>.

62 Replay is a candidate mechanism for this process, where a relative increase in prospective replay strength of one  
63 trajectory over another might relate to a bias towards deciding to approach versus avoid.

64 Here, we employed magnetoencephalography (MEG) to investigate whether there is an asymmetry between the  
65 replay of rewarding and aversive path sequences during planning. We designed a gambling-style task in which  
66 participants made decisions to either approach or avoid an uncertain environment containing paths leading to  
67 either gain or loss. By decoding rapid sequential replay related to sequences of transitioned states, we reveal a  
68 striking replay asymmetry that reflects prospective evaluations during planning and predicts trial-by-trial decision-  
69 making.

## 70 Results

### 71 Expected value guides decision-making

72 Participants learnt the structure of an environment containing two sequences (hereafter referred to as “paths”)  
73 containing three images (hereafter referred to as “states”), each with an associated integer value (**Extended Data**  
74 **Fig. 1C**). In a gambling-type scenario, where the overall task goal was to earn as many points as possible,  
75 participants could choose to either “approach” the environment, thereby probabilistically transitioning to one of  
76 the two paths, or “avoid” the environment entirely (receiving a guaranteed sum of 1 point).

77 To make a rational choice, participants needed to mentally simulate a prospective accumulation of points along  
78 each path. The total value of each path was dependent on a visual cue presented at the beginning of each trial,  
79 which also guided participants towards a sequential evaluation of each path in a forwards direction (**Fig. 1**; see  
80 description of “odd rule” in **Methods**). In a majority of trials, one of the two paths resulted in an overall gain and  
81 the other in an overall loss. The likelihood of transitioning to either of the two paths (conditional on participants  
82 choosing to approach) spanned five probabilities (10-90%, 30-70%, 50-50%, 70-30%, and 90-10%), and these were  
83 displayed on screen during an allowed 30-second planning period.

84 If participants chose to approach, a screen then displayed which of the two available paths had been selected, as  
85 determined by the path transition probabilities displayed during the planning phase (**Fig. 1A**). Participants then  
86 deterministically transitioned to each state along the selected path, with the state value and cumulative sum of  
87 points along the trajectory displayed on-screen. Note that the first four trials of each block were forced-choice to  
88 approach, serving as a reminder of the images representing each state (images were replaced by text labels in all  
89 other free-choice trials to control for visual exposure) and their associated integer values (the value of one state  
90 from each of the two paths was updated at the beginning of each block). If participants chose to avoid, a screen  
91 was then displayed indicating that participants had earned one point.

92 In the task, rational decision-making required calculating the expected value of approach (i.e., the sum of points  
93 for paths 1 and 2, weighted by their probabilities) and then choosing to approach only if the overall expected value  
94 is greater than a certain value granted after choosing to avoid (i.e.,  $\geq 1$ ). We calculated the accuracy of participants'  
95 choices by comparing them to perfectly rational choice behaviour. Simulations of different behavioural strategies  
96 showed that learning path values from experience, as opposed to the more cognitively burdensome optimal  
97 strategy of sequentially summing state values, could achieve approximately 69% to 85% accuracy (see  
98 **Supplementary methods** for model simulations). Moreover, only considering one, but not both, paths when  
99 computing the expected value achieved an hypothetical mean accuracy of approximately 63% (range = 53% to  
100 75%). Two of 26 participants performed at 47.55% and 51.37% accuracy, respectively, and thus were excluded  
101 from all subsequent analyses, except for evaluation of replay for an overall state map.

102 We expected an ordered reactivation of state transitions to reflect the repeated visual experience of paths in  
103 sequential order (i.e., during learning, as well as the walkthrough phases of decision trials), as shown by previous  
104 studies<sup>17,47</sup>, as opposed to reflecting a conscious mental calculation performed during planning. The sequential  
105 nature of determining path value was therefore a design feature that served to encourage perception of temporal  
106 order in the relationships between states, as well as provide a sufficient level of task difficulty. Moreover, the  
107 chosen design aligned with previous work using paradigms that incorporate a cumulative sum calculation<sup>12</sup>, as  
108 well as investigations of spatial replay (which is inherently sequential<sup>1-7</sup>).

109 Participants performed significantly above chance, with 76.07% accuracy on average (SD = 7.35%, range = 60.27%  
110 to 89.73%;  $t(23) = 17.373$ ,  $p < 0.001$ ; **Fig. 2A**), correctly approaching when the expected value was 2.386 on  
111 average (SD = 0.57) and avoiding when the expected value was -1.552 on average (SD = 0.556;  $t(23) = 21.152$ ,  $p$   
112  $< 0.001$ ; **Fig. 2B**). Overall, participants tended to approach more (57.15% of trials) than avoid (42.85%;  $t(23) =$   
113  $4.176$ ,  $p < 0.001$ ), consistent with reward-seeking or information-seeking behaviour. Experimental protocols were  
114 designed so that the expected value of approaching was  $> 1$  on 50% of trials (**Extended Data Fig.1E**). As such,  
115 accuracy was significantly lower on trials where participants chose to approach (74.59%, SD = 7.53%) than avoid  
116 (79.27%, SD = 8.07%;  $t(23) = -3.190$ ,  $p = 0.004$ ; **Fig. 2A**). Participants were also significantly faster in their decision  
117 to approach (M = 8.446 seconds, SD = 2.03) than to avoid (M = 8.975 seconds, SD = 2.424;  $t(23) = -2.319$ ,  $p =$   
118  $0.030$ ; **Fig. 2C**).

119 In the experimental design, there was consistency as to which of the two paths culminated in a reward or loss.  
120 Hence, for the first half of the experiment, path 1 resulted in reward and path 2 resulted in loss, and vice versa for  
121 the second half of the experiment (two protocols were used and counterbalanced across participants; see  
122 **Methods**). To encourage active engagement in sequential planning, rather than merely learning this tendency, we  
123 included catch trials (5%) where both paths either led to a reward or to a loss. Behavioural modelling of different  
124 strategies was consistent with participants performing online calculations (winning mental arithmetic model: N =  
125 15/24) as opposed to a strategy of caching learned values (N = 9/24; see **Supplementary methods**). Thus, for a

126 majority of participants (N = 20/24), choice behaviour was best explained by a model in which both paths were  
127 considered when computing expected value, while that of the remaining participants was best explained by  
128 models that either reflected a consideration of only path 1 (N = 2/24) or only the path perceived to be consistently  
129 punishing (N = 2/24). Thus, a majority of participants engaged in sequential planning by mentally accumulating  
130 points along each path, with the majority considering both paths (rather than just one path, in an effort to conserve  
131 cognitive resources) when deliberating. Note that, of the two participants excluded from path-specific analyses  
132 due to overall poor task accuracy, one was best explained by a null model (i.e., a general bias towards  
133 approaching, irrespective of expected outcome) and the other by a caching strategy that considered aversive  
134 paths alone.

135 We next constructed a multilevel logistic regression model to more precisely examine how path values and  
136 transition probabilities influenced trial-by-trial decision-making. In this model, trial-by-trial choice was predicted  
137 by a three-way interaction between the value of the path with the highest prospective value (i.e., the rewarding  
138 path), the value of the path with the lowest prospective value (i.e., the loss path), and the probability of  
139 transitioning to the rewarding path (the probability of transitioning to one path was always relative to the other).  
140 We also included response time (RT) as a fixed effect, as well as certainty of path transition probabilities on each  
141 trial (uncertain: 50-50%, moderately certain: 30-70% or 70-30%, very certain: 10-90% or 90-10%).

142 Approach choices were significantly predicted both by the probability of transitioning to a more rewarding path  
143 ( $\beta = 6.460$ ,  $p < 0.001$ ) and by larger prospective rewards ( $\beta = 0.113$ ,  $p < 0.001$ ; **Fig. 2D**). Thus, participants  
144 approached environments containing larger rewards more when the probability of transitioning to reward was  
145 higher ( $\beta = 4.991$ ,  $p < 0.001$ ). Although the magnitude of potential loss also predicted decision-making  
146 (participants were more likely to choose to avoid when potential losses were larger:  $\beta = 0.053$ ,  $p = 0.011$ ), there  
147 was no interaction with transition probability ( $\beta = -0.028$ ,  $p = 0.744$ ; **Fig. 2E**). These findings support the idea that  
148 decision-making was guided by the total value of reward and loss paths, as well as the probability of transitioning  
149 to a rewarding path. We also observed a significant effect of certainty ( $\beta = 0.317$ ,  $p < 0.001$ ), such that participants  
150 were more likely to approach when transition probabilities were more certain overall (i.e., 90-10% or 10-90%, as  
151 opposed to 50-50%).

152 Finally, given that participants were more likely to approach when rewarding paths were more probable, we also  
153 tested whether participants experienced rewarding paths more frequently than aversive paths. On average,  
154 participants transitioned to a rewarding path 107 times (SD = 11) and to an aversive path 23 times (SD = 8), a  
155 difference that was significant ( $t(23) = 25.577$ ,  $p < 0.001$ ). Importantly, due to our counterbalanced design there  
156 was no significant difference in the likelihood of experiencing path 1 (M = 46, SD = 8) or path 2 (M = 46, SD = 7;  
157  $t(23) = -0.046$ ,  $p = 0.964$ ).

## 158 Forward replay during planning

159 Our primary research questions with regard to replay were: 1) whether there is a sequential reactivation of state-  
160 to-state path transitions during planning, 2) whether this is influenced by each path's perceived value, and 3)  
161 whether this, in turn, relates to a subsequent choice to approach or avoid. In an initial functional localiser task, we  
162 trained classifiers on visually-evoked response fields (measured using MEG) for six unique state images (see **Fig.**  
163 **3A-C** and **Extended Data Fig. 2**). Importantly, these state neural signatures were captured prior to participants  
164 learning the order of states in each sequence. Next, we applied each state classifier to MEG data acquired during  
165 the planning period of each decision trial, producing time series of decoded state reactivation (**Fig. 3D**). Using  
166 general linear modelling, we assessed evidence for temporally-ordered reactivation of each state pair (A-B and B-  
167 C in path 1, and D-E, and E-F in path 2) across different time intervals (10 to 600 ms, in steps of 10 ms), in both a  
168 forwards and backwards direction. We refer to this as "sequenceness", our index of replay.

169 As a first step, we asked whether there was evidence for replay of the entire state space (i.e., average  
170 sequenceness across all four transitions), discarding the first four trials in each block as these were forced-  
171 choice. We observed maximal forward state-to-state reactivation at 60 ms intervals (or "lags"), and maximal  
172 backward state-to-state reactivation at 110 ms (**Fig. 3E**). We then computed a forward-minus-backward  
173 sequenceness measure to remove common noise and increase sensitivity. A significance threshold generated by  
174 random permutations (see **Methods**) provided evidence for significant forward replay at 20 to 90 ms state-to-  
175 state intervals, indicating the state space was replayed during planning with a rapid temporal compression akin  
176 to that reported in previous studies<sup>12,14,15,17,47</sup>. Notably, we did not observe significant forwards replay at longer  
177 state-to-state intervals of up to 3 seconds, where this might be more indicative of conscious memory retrieval  
178 processes during path evaluation or choice deliberation (**Extended Data Fig. 3**).

## 179 Replay is modulated by prospective reward and loss

180 Having found evidence for forwards replay during planning, we next asked whether we could differentiate replay  
181 for paths that culminated in either a reward or a loss. For each trial, we averaged sequenceness across the two  
182 transitions present within each path (**Fig. 4A**). We then entered these trial-by-trial estimates of path replay at the  
183 significant state-to-state intervals identified within our previous analysis (20 to 90 ms) into a series of linear  
184 mixed-effects models that accounted for effects of subject, replay interval, and trial duration (i.e., response time;  
185 see **Supplementary methods** for detailed model specification).

186 We first asked whether the expression of replay was influenced by an eventual choice to approach or avoid.  
187 Overall, rewarding paths were replayed more strongly than aversive paths during planning ( $\beta = 0.014$ ,  $p < 0.001$ ;  
188 **Fig. 4B**). Notably, there was a significant interaction with choice ( $\beta = -0.018$ ,  $p < 0.001$ ) showing this was  
189 particularly the case when participants made an eventual decision to avoid ( $EMM_{\Delta} = -0.008$ ,  $p < 0.001$ ). Replay

190 strength did not differ between reward and loss paths when participants planned to approach ( $EMM_{\Delta} = 0.002$ ,  $p$   
191  $= 0.362$ ). Thus, replay preceding a choice to avoid was stronger for potential paths leading to reward than for  
192 potential paths leading to punishment.

193 We next asked whether replay was modulated by factors other than choice; namely, recent path experience or the  
194 probability of transitioning to either path irrespective of path value. We operationalised recent experience as the  
195 number of trials within a block since participants last visited a particular path. We constructed a model in which  
196 replay was predicted by path experience (log-transformed to address a positive skew), path type (reward or loss),  
197 and path transition probability. Intriguingly, we found an interaction between path type and experience ( $\beta = 0.010$ ,  
198  $p < 0.001$ ), showing that rewarding paths were more strongly replayed when they had been less recently  
199 experienced, whereas loss paths were more weakly replayed (**Extended Data Fig. 4B**). Path transition probability  
200 modulated this effect ( $\beta = 0.014$ ,  $p = 0.010$ ), such that less recently experienced rewarding paths were even more  
201 strongly replayed when the upcoming transition probability was higher.

202 In our next model, we assessed whether path replay, irrespective of reward or loss, was modulated by its transition  
203 probability. We modelled replay of each path per trial as being predicted by its transition probability, as well as the  
204 subsequent choice made on each trial. We found no evidence for an effect of path transition probability on replay  
205 ( $\beta < 0.001$ ,  $p = 0.830$ ), regardless of which choice was being planned ( $\beta = -0.004$ ,  $p = 0.190$ ; **Extended Data Fig.**  
206 **4C**). This indicates that participants' beliefs about which path was more likely to be experienced did not impact  
207 the strength of replay.

208 Lastly, we asked whether evidence for a conscious retrieval of states during choice deliberation influenced the  
209 strength of path replay. Although our behavioural strategy modelling suggested participants did not have a bias  
210 towards evaluating one path more than another (**Extended Data Fig. 5**), we speculated that participants might  
211 differentially recollect states belonging to rewarding or aversive paths *after* appraising each path's value, as a way  
212 of simulating future outcomes during choice deliberation. We computed a measure of overall state reactivation  
213 throughout planning as an indicator of memory reactivation that might, in principle, be akin to conscious memory  
214 retrieval. We found that, overall, states belonging to paths leading to reward were reactivated more strongly overall  
215 ( $\beta < 0.001$ ,  $p = 0.027$ ; **Extended Data Fig. 6A**), but, crucially, a significant effect of choice and path type on replay  
216 remained even after accounting for such overall state reactivation ( $\beta = -0.008$ ,  $p < 0.001$ ; **Extended Data Fig. 6B**).

## 217 Replay predicts approach and avoidance

218 Stronger replay for rewarding paths when subjects planned to avoid indicates a relationship between the content  
219 of replay and subsequent decision-making. To investigate this further, we computed a measure of "differential"  
220 replay that captures a difference in the expression of sequenceness between each prospective path on a trial-by-  
221 trial basis. Specifically, we subtracted loss path replay from reward path replay, such that more positive

222 differential replay indicates a bias towards replaying paths leading to reward, and vice versa for more negative  
223 differential replay.

224 Using this differential replay measure, we modelled how replay content changed conditional on the choice being  
225 planned (i.e., to approach or avoid), as well as the environment prospects (i.e., the cumulative gain or loss for  
226 each path, and the probability of transitioning to each path). To simplify the model, we used the expected value  
227 of approaching on each trial as a summary metric of an environment's prospects (equivalent to the total sum of  
228 points for each path weighted by their respective transition probabilities). To then predict trial-by-trial decision-  
229 making, we constructed a model that allowed expected value to interact with differential replay at all significant  
230 replay intervals (20 to 90 ms), as well as the certainty of path transition probabilities and response times.

231 At a behavioural level, we observed a sigmoidal relationship between expected value and choice ( $\beta = 0.432$ ,  $p <$   
232  $0.001$ ; **Fig. 4C**), such that participants were more likely to approach when the associated expected value was  $\geq$  -  
233 1.2. This is below a rational indifference point of 1, indicating participants were more likely to approach  
234 environments with poorer prospects overall. Additionally, participants were more likely to approach when path  
235 transition probabilities were more certain ( $\beta = 0.270$ ,  $p < 0.001$ ).

236 At a neural level, trial-by-trial differential neural replay predicted choice ( $\beta = -0.713$ ,  $p < 0.001$ ), such that  
237 participants were more likely to approach when differential replay during planning was less positive, reflecting a  
238 bias towards replaying paths leading to potential loss and/or a bias away from replaying paths leading to potential  
239 reward. Importantly, this effect of differential replay on decision-making interacted with expected value ( $\beta = 0.133$ ,  
240  $p = 0.008$ ), such that a bias away from replaying paths leading to reward was even more pronounced when  
241 participants planned to approach environments with a more negative expected value.

242 Our use of a difference measure precludes knowing whether the above effect was driven by diminished replay of  
243 reward paths or enhanced replay of loss paths. To unpack this, we duplicated our model but replaced differential  
244 replay with two separate predictors, one for reward path replay and one for loss path replay, with each separately  
245 interacting with expected value. This revealed that path replay for reward and loss had opposing interactions with  
246 expected value, such that planning to approach a more hazardous environment (i.e., negative expected value)  
247 was predicted by enhanced replay of paths leading to loss ( $\beta = 0.120$ ,  $p = 0.090$ ) and an attenuated replay of paths  
248 leading to reward ( $\beta = -0.146$ ,  $p = 0.031$ ; **Fig. 4C**). Moreover, as highlighted by our earlier analyses of replay and  
249 path value, when participants planned to approach, replay of reward paths was significantly reduced ( $\beta = -1.232$ ,  
250  $p < 0.001$ ). Replay of loss paths did not predict decision-making ( $\beta = 0.189$ ,  $p = 0.313$ ). Thus, the content of replay  
251 predicted subsequent decisions such that when, participants exhibited more rational decision-making (i.e.,  
252 choosing to avoid when the expected value of an approach choice was lower), paths leading to reward were  
253 selectively replayed. By contrast, reduced replay of reward paths and relatively stronger replay of loss paths was  
254 associated with participants being more likely to approach riskier environments.



## 255 Trait anxiety and risk aversion

256 Next, we tested an hypothesis that a relationship between differential replay during planning and deciding to  
257 approach a risky environment would be amplified in participants with higher trait anxiety and/or a higher  
258 propensity towards risk-taking. An independent components analysis on subjects' self-report questionnaires  
259 yielded one component representing anxiety and another representing risk-aversion (see **Methods**). Based upon  
260 this, we then constructed a model in which these personality traits were allowed to interact with both differential  
261 replay during planning and expected value to predict future decision-making. We again included the degree of  
262 certainty about the path transition probabilities in the model, as well as trial duration (i.e., response time).

263 Within this model, anxiety and risk-aversion alone did not predict decision-making ( $\beta = -0.028$ ,  $p = 0.644$  and  $\beta =$   
264  $0.073$ ,  $p = 0.284$ , respectively), although there was a significant increase in approach rate at higher expected  
265 values (indicating more conservative decision-making) in participants with higher risk aversion ( $\beta = 0.008$ ,  $p =$   
266  $0.048$ ). Instead, both anxiety and risk aversion significantly modulated the relationship between differential replay  
267 and decision-making. More anxious ( $\beta = -0.314$ ,  $p = 0.003$ ; **Fig. 4D**) and more risk-averse participants ( $\beta = -0.377$ ,  
268  $p < 0.001$ ; **Fig. 4E**) showed a greater likelihood of approaching when replay was biased away from rewarding paths  
269 ( $\beta = -0.314$ ,  $p = 0.003$ ). For more risk-averse participants, this was the case regardless of expected value ( $\beta = -$   
270  $0.018$ ,  $p = 0.540$ ), whereas for more anxious participants this was predominantly the case when expected value  
271 was lower ( $\beta = 0.096$ ,  $p = 0.014$ ).

272 We repeated this model using separate interacting predictors for reward and loss path replay to detail what was  
273 driving the above effects. The model revealed that replay for paths leading to loss ( $\beta = 0.450$ ,  $p = 0.003$ ), but not  
274 reward ( $\beta = -0.199$ ,  $p = 0.189$ ), was boosted in more anxious participants when approaching more aversive  
275 environments. Similarly, replay for paths leading to loss ( $\beta = 0.483$ ,  $p < 0.001$ ), but not reward ( $\beta = -0.225$ ,  $p =$   
276  $0.084$ ), was boosted for more risk-averse participants when planning to approach any environment. Additionally,  
277 more risk-averse participants had diminished replay of rewarding paths when planning to approach more aversive  
278 environments, while more risk-seeking participants had diminished replay of rewarding paths when planning to  
279 approach more lucrative environments ( $\beta = -0.132$ ,  $p = 0.004$ ). In contrast, more anxious participants had stronger  
280 replay of loss paths when planning to approach more aversive environments ( $\beta = -0.201$ ,  $p < 0.001$ ). This suggests  
281 that the more negative differential replay in participants with higher trait anxiety during planning was driven by an  
282 increase in loss path replay rather than a decrease in reward path replay, while the opposite was true for approach  
283 planning in more risk-averse participants.

## 284 Fronto-temporal theta activity underlies replay during planning

285 In a final analysis, we estimated the spatial sources of activity underlying onset of replay events. We defined a  
286 replay "event" as an above-chance reactivation of one state followed by reactivation of the following state within

287 a 20 to 90 ms lag, with additional stringent criteria (see **Methods**). We reconstructed source activity in either the  
288 theta (4 to 8 Hz) or high gamma (120 to 150 Hz) frequency band based upon a priori evidence for expression of  
289 hippocampal theta-related replay during planning<sup>48,49</sup>, as well as high-frequency sharp-wave ripple events in  
290 hippocampus related to planning<sup>19</sup>.

291 Across the whole brain ( $p < 0.05$ , FWE-corrected), there was a significant increase in theta power in the right  
292 thalamus (peak MNI: 5, -26, 8), as well as a cluster spanning the left middle temporal gyrus that overlapped left  
293 posterior hippocampus (peak MNI: -40, -31, -2). We also observed significant theta activity in dorsolateral  
294 prefrontal cortex (DLPFC; peak MNI: -35, 29, 28), right anterior cingulate cortex (ACC; peak MNI: 10, 49, 13),  
295 striatum (peak MNI: -15, 4, 13), and inferior occipital cortex (0, -101, -12). In contrast, we did not observe significant  
296 source activity in the gamma frequency range during replay events. We also investigated whether theta or high-  
297 gamma activity during replay events covaried with each subject's trait anxiety or overall performance accuracy,  
298 but we did not observe any significant effects.

299 The increased theta activity in medial temporal lobe accords with studies in rodent hippocampus, where a rapid  
300 "look-ahead" of spatial trajectories during route planning is reflected by rapid hippocampal replay events bounded  
301 by theta cycles<sup>49</sup>. In humans, theta activity in hippocampus and medial temporal lobe has been observed during  
302 prospective replay events when participants plan to avoid aversive outcomes<sup>48</sup>, similar to the present study. Other  
303 studies in humans have localised replay onset during post-task rest periods (associated with memory  
304 consolidation of a cognitive map) to left hippocampus in the gamma frequency band<sup>14,15,47</sup>. This pattern is in line  
305 with the notion that planning-related replay in medial temporal cortex during is subserved by theta activity,  
306 whereas replay related to memory consolidation at rest is more closely linked to high-frequency sharp wave ripple  
307 events<sup>18,49-51</sup>.

308 Our results support previous evidence for a role for ACC<sup>52</sup>, DLPFC<sup>53,54</sup>, striatum<sup>55</sup>, and inferior occipital cortex<sup>15</sup> in  
309 prospective replay during planning that involves elements such as rule-switching or reward re-evaluation.  
310 Intriguingly, our results also hint at the thalamus being a significant source of replay-related theta activity. The  
311 thalamus purportedly coordinates reward-guided decision-making processes across hippocampus, medial  
312 temporal lobe, and prefrontal cortex<sup>56</sup>, and thus might reasonably be involved in long-range communication of  
313 ordered state reactivation across these areas during planning.

## 314 Discussion

315 In rodent studies, replay content during planning has been found to reflect paths that should be pursued<sup>23,24</sup> as  
316 well as those that should be avoided<sup>33</sup>. Here, in the context of an approach-avoidance conflict in humans, we find  
317 that the content of forward replay during planning flexibly predicted subsequent decisions. Participants were  
318 more likely to avoid when replay was relatively stronger for paths leading to reward, and more likely to approach

319 when replay was relatively stronger for paths leading to loss, an effect most pronounced for risky environments  
320 (i.e., there was a negative expected value of approaching). Our findings indicate a role for replay during planning  
321 under uncertainty, where the relative strength of replay for paths leading to reward and loss is weighted towards  
322 counterfactual outcomes relating to a current plan to approach or avoid.

323 Based on rodent studies, we had expected prospective replay content to reflect the goals of approach (to obtain  
324 reward) and avoidance (to avoid punishment), such that replay would increase for rewarding paths being  
325 pursued<sup>21,23,24,26,50</sup> and for punishing paths being avoided<sup>33</sup>. Instead, we observed the opposite pattern, albeit in an  
326 environment that contained both reward and loss paths. Preceding a decision to avoid, replay was increased for  
327 paths that would lead to a foregone reward. By contrast, replay of paths leading to prospective reward was  
328 decreased preceding a decision to approach. Indeed, when there was greater risk associated with approach (i.e.,  
329 a negative expected value), replay increased for paths leading to potential loss. This suggests a relationship  
330 between the content of prospective replay and rational decision-making under risk, where boosted replay of  
331 rewarding paths predicted rational avoidance and replay of more punishing paths predicted irrational approach.

332 Our findings echo a recent theoretical account which proposed replay provides a pessimistic reminder of  
333 counterfactual outcomes to a model-free learning system<sup>37</sup>. This proposal finds support in observations of  
334 increased replay for paths previously – but not currently – rewarded<sup>35,36</sup>. Furthermore, behavioural modelling has  
335 linked replay to model-based planning, such that replay of sub-optimal outcomes of a given choice (i.e.,  
336 “pessimistic” replay) promotes more rational model-free decision-making by ensuring that negative outcomes of  
337 unchosen actions are not forgotten<sup>22,37,48</sup>. Intriguingly, we also observed enhanced replay for paths leading to  
338 counterfactual outcomes (albeit the hypothetical outcomes of a planned decision, rather than the observed  
339 outcomes of a previous decision), though the design of our experiment does not allow us to draw conclusions  
340 regarding a contribution of  $r$  to either model-based versus model-free learning mechanisms. Our task entailed a  
341 high degree of cross-trial volatility in state values and transition probabilities that rendered model-free learning of  
342 state-action contingencies futile, as participants needed to adopt a model-based strategy that explicitly  
343 considered path values and transition probabilities. Moreover, we found mixed evidence for whether paths were  
344 more strongly replayed when they were more susceptible to being forgotten. While we found paths were replayed  
345 more strongly when they had not been recently experienced, this was the case solely for rewarding, but not  
346 aversive, paths. Additionally, under an assumption that replay prevents forgetting of sequences and their  
347 associated values, then planning to avoid should theoretically increase replay of both paths (as neither will be  
348 experienced) while planning to approach should increase replay of the less probable path. However, this was only  
349 the case for rewarding paths, suggesting that counterfactual replay during model-based planning is not  
350 adequately explained by a role in memory maintenance<sup>14,15,18,29–31,57</sup>.

351 An alternative explanation for the pattern of replay we observed is that it reflects an anxiety-related simulation of  
352 counterfactual outcomes during planning. This would explain counterfactual replay being associated with

353 irrational decisions to approach under riskier conditions (i.e., when the expected value of approaching was  
354 negative), an effect most pronounced in participants with higher self-reported trait anxiety. Similarly, for  
355 participants who self-reported higher trait risk-aversion replay was biased towards paths leading to loss when  
356 planning to approach, regardless of the expected value of approaching. Replay has been speculated to play a role  
357 in clinical anxiety and depression<sup>41</sup> and our study provides tentative evidence for a relationship between  
358 differential replay and both trait anxiety and risk-aversion.

359 Dispositional anxiety is associated with a heightened, and sometimes uncontrollable, simulation of potential past  
360 (rumination) or future (worry) aversive events<sup>40,58</sup>. A functional role for replay in selectively sampling a prospective  
361 environment during planning provides a plausible explanation for a bias towards more aversive outcomes in  
362 people who have a greater tendency to worry<sup>40</sup>. Indeed, people with higher social anxiety engage in  
363 “counterfactual” updating, entailing greater deliberation of outcomes that have not, or will not, be experienced<sup>59</sup>.  
364 Thus, a simulation of “what if ” scenarios maps closely with our finding that replay content reflects a worst-case  
365 scenario of a plan to approach (i.e., the possibility of being punished) or avoid (i.e., foregoing potential reward).  
366 Note that more anxious and more risk-averse participants did not make more erroneous approach decisions  
367 overall, and an effect of anxiety and risk-aversion was only discernible at the neural level. Thus, our findings do  
368 not provide support for a suggestion that counterfactual replay drives a change in policy per se. Moreover, as our  
369 sample consisted solely of healthy controls, future studies involving participants with anxiety disorders, who show  
370 irrational risky choice behaviour<sup>60</sup>, could determine the extent to which replay relates to anxiety-modulated model-  
371 based decision-making.

372 An important caveat to our study is that reward may have been perceived by participants as more salient than  
373 loss, in line with participants’ choices being more sensitive to probability and magnitude of reward than that of  
374 loss. Playing to accumulate monetary rewards, as opposed to avoiding monetary losses, has been shown to  
375 enhance the utility of reward<sup>61</sup>. This might explain why replay reflected a worst-case scenario of choosing to avoid  
376 (i.e., foregoing potential reward) across all trials, irrespective of expected value. By contrast, replay reflected the  
377 worst-case scenario of choosing to approach (i.e., transitioning to a loss path) only when the expected value of  
378 approaching was more negative. An emphasis on reward might also explain why a relationship between path  
379 replay and memory maintenance was more evident for rewarding paths (as discussed above) but not punishing  
380 paths, in line with other recent findings<sup>36</sup>. Employing a variant of the current design using more arousing positive  
381 and negative stimuli (e.g., electric shocks or affective visual stimuli) could help adjudicate between these  
382 possibilities.

383 Overall, we present novel evidence for a relationship between the expression of replay and decision-making under  
384 uncertainty. A path prioritisation in prospective replay reflected a worst-case scenario of a decision to approach  
385 (increased replay for loss paths) or avoid (increased replay for reward paths). Our findings align with recent  
386 observations that replay reflects counterfactual outcomes associated with prospective decision-making<sup>35-37</sup> and

387 extends this to a domain in which choices to pursue reward also carry a risk of punishment. Scenarios such as  
388 this are particularly pertinent to survival where an outcome might be critical for the viability of an agent<sup>62</sup>, as well  
389 as to understanding anxiety-related disorders that are characterised by an over-simulation of improbable, but  
390 often catastrophic, events<sup>40</sup>.

## 391 Data availability

392 Data are freely available on the Open Science Framework: <https://osf.io/6ndu9/>.

## 393 Code availability

394 All code for the experimental paradigm and analysis pipeline is freely available on GitHub:  
395 <https://github.com/jjmcadyen/approach-avoid-replay>.

## 396 Acknowledgements

397 The authors thank Toby Wise and Paul Sharp for their helpful discussions about the study design. This work is  
398 supported by the Wellcome Trust (098362/A/12/Z and 091593/Z/10/Z supporting RJD and JM, respectively). YL  
399 is supported by National Science and Technology Innovation 2030 Major Program (2022ZD0205500) and National  
400 Natural Science Foundation of China (32271093). The Max Planck University College London Centre for  
401 Computational Psychiatry and Ageing Research is a joint initiative supported by University College London and  
402 the Max Planck Society. The Wellcome Centre for Human Neuroimaging is supported by core funding from the  
403 Wellcome Trust (203147/Z/16/Z). The funders had no role in study design, data collection and analysis, decision  
404 to publish, or preparation of the manuscript.

## 405 Author contributions

406 JM designed the experiment with input from YL. JM collected the data, and JM and YL wrote the analysis code.  
407 JM and YL interpreted the data with input from RJD. JM wrote the manuscript with input and edits from YL and  
408 RJD.

## 409 Competing interests

410 The authors declare no competing interests.

## 411 Figure legends

412 **Figure 1. Decision trials. (A)** Participants began each trial in a planning phase. Here, they used the presented  
413 information (the odd rule states and the path transition probabilities) to mentally calculate the total outcome for  
414 each path and evaluate the utility of an approach vs avoid decision. This calculation involves summing the value  
415 ( $v$ ) of each state ( $s$ ) across each path, taking into account the 'odd rule', and multiplying the final sum ( $R_{path}$ ) by  
416 the path transition probabilities ( $P_{path}$ ), as described in B. The order of images and their respective values were  
417 learned during an initial training phase (Extended Data Fig. 1). MEG data from this planning period provided the  
418 focus for our replay analysis. If participants chose to approach, a screen then appeared displaying which of two  
419 potential paths they had probabilistically transitioned to ("Transition" screen), and participants then observed an  
420 animation of this sequence ("Walkthrough" screens). During this walkthrough, the number of points gained or lost  
421 at each state (light blue numbers), as well as the cumulative sum of points up to and including each state (dark  
422 blue numbers), was shown below the state image. Note that images were only shown in forced-choice trials, while  
423 text labels were shown in all other trials. The final sum of points for the sequence was then shown ("Outcome"  
424 screen). If participants chose to avoid, a fixed increase of one point was shown ("Safe outcome" screen). **(B)** The  
425 "odd rule" was introduced to reinforce the temporal order relationships between states by having participants  
426 appraise each sequence in a forwards direction. The rule was always applied to one state from each path, and  
427 this was indicated to participants on-screen during planning. The odd rule entailed that, if the cumulative sum of  
428 points collected up to (and including) a particular state was an odd number, then the sign of the sum would then  
429 be reversed (i.e., multiplied by -1) with this sum being carried over to any subsequent states in that path. Thus,  
430 the odd rule could significantly alter the total number of points collected along each path, depending on which  
431 state the odd rule was applied to, and enforced a need for online calculation. For example, using the values of  
432 path 2 illustrated in A, applying the odd rule to state 2 results in -3 points ( $s_1: 0 + 4 = 4 \rightarrow s_2: 4 + 1 = 5$ , which is  
433 odd and so the sign is reversed to give  $-5 \rightarrow s_3: -5 + 2 = -3$ ), whereas applying the odd rule to state 1 results in 7  
434 points ( $s_1: 0 + 4$ , as the sign is not reversed  $\rightarrow s_2: 4 + 1 = 5 \rightarrow s_3: 5 + 2 = 7$ ). A rational planner first calculates the  
435 cumulative sum of points along each path (taking the odd rule into account), multiplies these by the respective  
436 path transition probabilities (which varied trial to trial), and then decides based on a comparison between the  
437 expected value of approaching ( $EV_{app}$ ) and the expected value of avoiding ( $EV_{av}$ ).

438 **Figure 2. Behavioural results. (A)** Accuracy is defined as the proportion of trials wherein participant responses  
439 matched an optimal response, based upon expected value. Overall, participants (individual markers;  $N = 26$ ) made  
440 significantly more accurate avoid decisions than approach decisions (two-tailed  $t(25) = 4.023$ ,  $p = 4.591E-4$ ).  
441 Horizontal line indicates median and box bounds indicate 25<sup>th</sup> and 75<sup>th</sup> quantile. **(B)** The expected value of  
442 approaching was significantly higher when participants ( $N = 26$ ) chose to approach than when participants chose  
443 to avoid (two-tailed  $t(25) = 12.250$ ,  $p = 4.614E-12$ ). **(C)** Participants ( $N = 26$ ) were significantly faster to approach  
444 than to avoid (two-tailed  $t(25) = -2.360$ ,  $p = 0.026$ ). Boxplots indicate median and 25<sup>th</sup> and 75<sup>th</sup> percentiles of

445 average participant response times. **(D)** Approach rate estimated by a behavioural multilevel model showing  
446 participants were more likely to approach if the probability of transitioning to a rewarding path was higher,  
447 especially when prospective reward values were greater (error bars indicate 95% confidence interval). **(E)**  
448 Similarly, participants were more likely to approach if potential loss was lower, irrespective of path transition  
449 probability (error bars indicate 95% confidence interval). \*  $p < .05$ , \*\*  $p < .01$ , \*\*\*  $p < 0.001$ .

450 **Figure 3. State classification and replay analysis.** **(A)** Before learning the order of images along each path,  
451 participants viewed each image in an initial functional localiser task. The visually-evoked event-related fields  
452 (measured using MEG) are displayed for each of the 12 images, or “states” (6 were randomly assigned to each  
453 participant), averaged across participants (shaded error indicates standard error of the mean). **(B)** Using  
454 functional localiser MEG data, we created classifiers for each state, per participant (example participant shown).  
455 A classifier was a set of beta weights per sensor. **(C)** Using K-folds cross-validation, we assessed average  
456 accuracy of state classifiers per participant. Classifiers trained at a 120 ms time point produced the highest  
457 average accuracy overall (error bars indicate standard error of the mean). **(D)** Classifiers trained on either 110,  
458 120, or 130 ms (accounting for inter-subject variability in classifier performance) were applied to MEG data  
459 collected throughout the planning period of decision trials, producing matrices of predicted state reactivation per  
460 trial (example shown). **(E)** Using a two-level GLM approach, we estimated the intervals (or “lags”) between  
461 maximal reactivation of each state during planning, in a forwards (left) and backwards (middle) direction. Plots  
462 display the sequenceness estimates averaged across all four transitions (shaded error indicates standard error  
463 of the mean). The significance threshold is indicated by an horizontal dashed line. Significant forwards-minus-  
464 backwards replay occurred at state-to-state intervals of 20 to 90 ms, peaking at 60 ms.

465 **Figure 4. Replay of prospective reward and loss paths.** **(A)** Replay strength for paths leading to either reward  
466 (green) or loss (red) during planning, split according to whether participants subsequently chose to approach  
467 (left) or avoid (right). Data is averaged across trials and participants. Significant replay intervals are highlighted  
468 by the yellow box. The difference between reward and loss replay is also shown (black). **(B)** Estimated marginal  
469 means produced by a mixed-effects model ( $N = 24$  participants) predicting replay strength by the total value of a  
470 path (reward in green, loss in red) and the choice subsequently made by participants (approach or avoid). Error  
471 bars indicate standard error, and significance is given by a two-tailed statistic using a Satterthwaite approximation  
472 ( $p = 3.452E-6$ ). \*  $p < .05$ , \*\*  $p < .01$ , \*\*\*  $p < 0.001$  **(C)** Approach rate (y axis) predicted by a mixed-effects model  
473 containing expected value (x axis) and differential replay. When differential replay was more negative (red,  
474 indicating relatively stronger replay of loss than reward paths), participants were more likely to approach  
475 environments with poorer prospects (i.e., negative expected value). A similar model using separate predictors for  
476 reward and loss path replay showed participants were more likely to approach on trials with a negative expected  
477 value when replay of rewarding paths was attenuated (green dashed) and when replay of loss paths was enhanced  
478 (red solid). The indifference point (i.e., the point at which approach rate should be 50%) is displayed for rational  
479 agent behaviour (vertical dashed line). **(D)** Same as **C**, except that participants’ trait anxiety and risk-aversion

480 scores were included in the model. The interaction between differential replay and expected value on choice was  
481 driven predominantly by more anxious participants (low/high split is for visualisation purposes only). **(E)** Same as  
482 **D**, except data has been split into low and high risk-aversion. More risk-averse participants were more likely to  
483 approach when differential replay was more negative, regardless of expected value.

484 **Figure 5. Beamforming analysis on replay onsets.** Sources underlying the onset of replay events for any state-to-  
485 state transition included the middle temporal gyrus, hippocampus, anterior cingulate cortex (ACC), and thalamus.  
486 Significant activity not pictured: striatum, dorsolateral prefrontal cortex (DLPFC), and inferior occipital cortex.  
487 Viewing coordinates: left and middle = MNI [-30, -30, 3], right = MNI [5, 47, 7]. Clusters are thresholded at  $p < 0.05$ ,  
488 whole brain FWE-corrected.

## 489 References

- 490 1. Skaggs, W. E. & McNaughton, B. L. Replay of neuronal firing sequences in rat hippocampus during sleep  
491 following spatial experience. *Science* **271**, 1870–1873 (1996).
- 492 2. Louie, K. & Wilson, M. A. Temporally structured replay of awake hippocampal ensemble activity during rapid  
493 eye movement sleep. *Neuron* **29**, 145–156 (2001).
- 494 3. Diba, K. & Buzsáki, G. Forward and reverse hippocampal place-cell sequences during ripples. *Nat. Neurosci.*  
495 **10**, 1241–1242 (2007).
- 496 4. Foster, D. J. & Wilson, M. A. Reverse replay of behavioural sequences in hippocampal place cells during the  
497 awake state. *Nature* **440**, 680–683 (2006).
- 498 5. Lee, A. K. & Wilson, M. A. Memory of sequential experience in the hippocampus during slow wave sleep.  
499 *Neuron* **36**, 1183–1194 (2002).
- 500 6. Wilson, M. A. & McNaughton, B. L. Reactivation of hippocampal ensemble memories during sleep. *Science*  
501 **265**, 676–679 (1994).
- 502 7. Schuck, N. W. & Niv, Y. Sequential replay of nonspatial task states in the human hippocampus. *Science* **364**,  
503 (2019).
- 504 8. Wu, X. & Foster, D. J. Hippocampal replay captures the unique topological structure of a novel environment.  
505 *J. Neurosci.* **34**, 6459–6469 (2014).
- 506 9. Davidson, T. J., Kloosterman, F. & Wilson, M. A. Hippocampal replay of extended experience. *Neuron* **63**,  
507 497–507 (2009).
- 508 10. Jadhav, S. P., Kemere, C., German, P. W. & Frank, L. M. Awake hippocampal sharp-wave ripples support  
509 spatial memory. *Science* **336**, 1454–1458 (2012).
- 510 11. Karlsson, M. P. & Frank, L. M. Awake replay of remote experiences in the hippocampus. *Nat. Neurosci.* **12**,  
511 913–918 (2009).
- 512 12. Kurth-Nelson, Z., Economides, M., Dolan, R. J. & Dayan, P. Fast Sequences of Non-spatial State



- 513 Representations in Humans. *Neuron* **91**, 194–204 (2016).
- 514 13. Friston, K. & Buzsáki, G. The functional anatomy of time: What and when in the brain. *Trends Cogn. Sci.* **20**,  
515 500–511 (2016).
- 516 14. Nour, M. M., Liu, Y., Arumham, A., Kurth-Nelson, Z. & Dolan, R. J. Impaired neural replay of inferred  
517 relationships in schizophrenia. *Cell* **184**, 4315–4328.e17 (2021).
- 518 15. Liu, Y., Dolan, R. J., Kurth-Nelson, Z. & Behrens, T. E. J. Human Replay Spontaneously Reorganizes  
519 Experience. *Cell* **178**, 640–652.e14 (2019).
- 520 16. Wallach, H. *et al.* Coordinated hippocampal-entorhinal replay as structural inference. in *Advances In Neural*  
521 *Information Processing Systems 32 (NIPS 2019)* (eds. Wallach, H. *et al.*) vol. 32 13 (Neural Information  
522 Processing Systems (NIPS), 2019).
- 523 17. Liu, Y., Mattar, M. G., Behrens, T. E. J., Daw, N. D. & Dolan, R. J. Experience replay is associated with efficient  
524 nonlocal learning. *Science* **372**, (2021).
- 525 18. Ólafsdóttir, H. F., Bush, D. & Barry, C. The Role of Hippocampal Replay in Memory and Planning. *Curr. Biol.*  
526 **28**, R37–R50 (2018).
- 527 19. Buzsáki, G. Hippocampal sharp wave-ripple: A cognitive biomarker for episodic memory and planning.  
528 *Hippocampus* (2015).
- 529 20. Ólafsdóttir, H. F., Carpenter, F. & Barry, C. Task Demands Predict a Dynamic Switch in the Content of Awake  
530 Hippocampal Replay. *Neuron* **96**, 925–935.e6 (2017).
- 531 21. Xu, H., Baracskay, P., O’Neill, J. & Csicsvari, J. Assembly Responses of Hippocampal CA1 Place Cells  
532 Predict Learned Behavior in Goal-Directed Spatial Tasks on the Radial Eight-Arm Maze. *Neuron* **101**, 119–  
533 132.e4 (2019).
- 534 22. Eldar, E., Lièvre, G., Dayan, P. & Dolan, R. J. The roles of online and offline replay in planning. *Elife* **9**, (2020).
- 535 23. Singer, A. C., Carr, M. F., Karlsson, M. P. & Frank, L. M. Hippocampal SWR activity predicts correct decisions  
536 during the initial learning of an alternation task. *Neuron* **77**, 1163–1173 (2013).
- 537 24. Pfeiffer, B. E. & Foster, D. J. Hippocampal place-cell sequences depict future paths to remembered goals.  
538 *Nature* **497**, 74–79 (2013).
- 539 25. Ambrose, R. E., Pfeiffer, B. E. & Foster, D. J. Reverse Replay of Hippocampal Place Cells Is Uniquely  
540 Modulated by Changing Reward. *Neuron* **91**, 1124–1136 (2016).
- 541 26. Zheng, C., Hwaun, E., Loza, C. A. & Colgin, L. L. Hippocampal place cell sequences differ during correct and  
542 error trials in a spatial memory task. *Nat. Commun.* **12**, 3373 (2021).
- 543 27. Papale, A. E., Zielinski, M. C., Frank, L. M., Jadhav, S. P. & Redish, A. D. Interplay between Hippocampal  
544 Sharp-Wave-Ripple Events and Vicarious Trial and Error Behaviors in Decision Making. *Neuron* **92**, 975–982  
545 (2016).
- 546 28. Igata, H., Ikegaya, Y. & Sasaki, T. Prioritized experience replays on a hippocampal predictive map for  
547 learning. *Proc. Natl. Acad. Sci. U. S. A.* **118**, (2021).
- 548 29. Gupta, A. S., van der Meer, M. A. A., Touretzky, D. S. & Redish, A. D. Hippocampal replay is not a simple

- 549 function of experience. *Neuron* **65**, 695–705 (2010).
- 550 30. Ólafsdóttir, H. F., Barry, C., Saleem, A. B., Hassabis, D. & Spiers, H. J. Hippocampal place cells construct  
551 reward related sequences through unexplored space. *Elife* **4**, e06063 (2015).
- 552 31. Dragoi, G. & Tonegawa, S. Preplay of future place cell sequences by hippocampal cellular assemblies.  
553 *Nature* **469**, 397–401 (2011).
- 554 32. Singer, A. C. & Frank, L. M. Rewarded outcomes enhance reactivation of experience in the hippocampus.  
555 *Neuron* **64**, 910–921 (2009).
- 556 33. Wu, C.-T., Haggerty, D., Kemere, C. & Ji, D. Hippocampal awake replay in fear memory retrieval. *Nat.*  
557 *Neurosci.* **20**, 571–580 (2017).
- 558 34. Mattar, M. G. & Daw, N. D. Prioritized memory access explains planning and hippocampal replay. *Nature*  
559 *Neuroscience* vol. 21 1609–1617 Preprint at <https://doi.org/10.1038/s41593-018-0232-z> (2018).
- 560 35. Carey, A. A., Tanaka, Y. & van der Meer, M. A. A. Reward revaluation biases hippocampal replay content  
561 away from the preferred outcome. *Nat. Neurosci.* **22**, 1450–1459 (2019).
- 562 36. Gillespie, A. K. *et al.* Hippocampal replay reflects specific past experiences rather than a plan for  
563 subsequent choice. *Neuron* (2021) doi:10.1016/j.neuron.2021.07.029.
- 564 37. Antonov, G., Gagne, C., Eldar, E. & Dayan, P. Optimism and Pessimism in Optimised Replay. *PLOS*  
565 *Computational Biology* **18**, e1009634 (2022).
- 566 38. Quartz, S. R. Reason, emotion and decision-making: risk and reward computation with feeling. *Trends Cogn.*  
567 *Sci.* **13**, 209–215 (2009).
- 568 39. Bach, D. R. *et al.* Human Hippocampus Arbitrates Approach-Avoidance Conflict. *Curr. Biol.* **24**, 1435 (2014).
- 569 40. Gagne, C., Dayan, P. & Bishop, S. J. When planning to survive goes wrong: predicting the future and  
570 replaying the past in anxiety and PTSD. *Current Opinion in Behavioral Sciences* **24**, 89–95 (2018).
- 571 41. Heller, A. S. & Bagot, R. C. Is Hippocampal Replay a Mechanism for Anxiety and Depression? *JAMA*  
572 *Psychiatry* **77**, 431–432 (2020).
- 573 42. Aupperle, R. L. & Paulus, M. P. Neural systems underlying approach and avoidance in anxiety disorders.  
574 *Dialogues Clin. Neurosci.* **12**, 517–531 (2010).
- 575 43. Loijen, A., Vrijssen, J. N., Egger, J. I. M., Becker, E. S. & Rinck, M. Biased approach-avoidance tendencies in  
576 psychopathology: A systematic review of their assessment and modification. *Clin. Psychol. Rev.* **77**, 101825  
577 (2020).
- 578 44. Loh, E. *et al.* Parsing the Role of the Hippocampus in Approach–Avoidance Conflict. *Cereb. Cortex* **27**, 201–  
579 215 (2016).
- 580 45. Schumacher, A. *et al.* Ventral Hippocampal CA1 and CA3 Differentially Mediate Learned Approach-  
581 Avoidance Conflict Processing. *Curr. Biol.* **28**, 1318–1324.e4 (2018).
- 582 46. Schumacher, A., Vlassov, E. & Ito, R. The ventral hippocampus, but not the dorsal hippocampus is critical for  
583 learned approach-avoidance decision making. *Hippocampus* **26**, 530–542 (2016).
- 584 47. Wimmer, G. E., Liu, Y., Vehar, N., Behrens, T. E. J. & Dolan, R. J. Episodic memory retrieval success is

- 585 associated with rapid replay of episode content. *Nat. Neurosci.* **23**, 1025–1033 (2020).
- 586 48. Wise, T., Liu, Y., Chowdhury, F. & Dolan, R. J. Model-based aversive learning in humans is supported by  
587 preferential task state reactivation. *Sci Adv* **7**, (2021).
- 588 49. Wikenheiser, A. M. & Redish, A. D. Hippocampal theta sequences reflect current goals. *Nat. Neurosci.* **18**,  
589 289–294 (2015).
- 590 50. Johnson, A. & Redish, A. D. Neural ensembles in CA3 transiently encode paths forward of the animal at a  
591 decision point. *J. Neurosci.* **27**, 12176–12189 (2007).
- 592 51. Gupta, A. S., van der Meer, M. A. A., Touretzky, D. S. & Redish, A. D. Segmentation of spatial experience by  
593 hippocampal theta sequences. *Nat. Neurosci.* **15**, 1032–1039 (2012).
- 594 52. Momennejad, I., Otto, A. R., Daw, N. D. & Norman, K. A. Offline replay supports planning in human  
595 reinforcement learning. *Elife* **7**, (2018).
- 596 53. Kaefer, K., Nardin, M., Blahna, K. & Csicsvari, J. Replay of Behavioral Sequences in the Medial Prefrontal  
597 Cortex during Rule Switching. *Neuron* **106**, 154–165.e6 (2020).
- 598 54. Berners-Lee, A., Wu, X. & Foster, D. J. Prefrontal cortical neurons are selective for non-local hippocampal  
599 representations during replay and behavior. *J. Neurosci.* (2021) doi:10.1523/JNEUROSCI.1158-20.2021.
- 600 55. Lansink, C. S., Goltstein, P. M., Lankelma, J. V., McNaughton, B. L. & Pennartz, C. M. A. Hippocampus leads  
601 ventral striatum in replay of place-reward information. *PLoS Biol.* **7**, e1000173 (2009).
- 602 56. Vertes, R. P. Interactions among the medial prefrontal cortex, hippocampus and midline thalamus in  
603 emotional and cognitive processing in the rat. *Neuroscience* **142**, 1–20 (2006).
- 604 57. Tolman, E. C. Cognitive maps in rats and men. *Psychol. Rev.* **55**, 189–208 (1948).
- 605 58. Hirsch, C. R. & Mathews, A. A cognitive model of pathological worry. *Behav. Res. Ther.* **50**, 636–646 (2012).
- 606 59. Hunter, L. E., Meer, E. A., Gillan, C. M., Hsu, M. & Daw, N. D. Increased and biased deliberation in social  
607 anxiety. *Nat Hum Behav* (2021) doi:10.1038/s41562-021-01180-y.
- 608 60. Hartley, C. A. & Phelps, E. A. Anxiety and decision-making. *Biol. Psychiatry* **72**, 113–118 (2012).
- 609 61. Kuhnen, C. M. Asymmetric learning from financial information. *J. Finance* **70**, 2029–2062 (2015).
- 610 62. Mobbs, D., Headley, D. B., Ding, W. & Dayan, P. Space, Time, and Fear: Survival Computations along  
611 Defensive Circuits. *Trends Cogn. Sci.* **24**, 228–241 (2020).

## 612 Methods

### 613 Participants

614 The study was approved by the University College London Research Ethics Committee (9929/002). We recruited  
615 32 healthy volunteers via online advertisements to participate in the first session, which served as an opportunity  
616 to practice and as a screening point to exclude participants who found the memorisation or arithmetic in the task

617 too difficult (see **Methods, Experimental task**). We excluded 1 participant who scored < 80% accuracy when tested  
618 on the image order, and 4 participants who scored < 60% accuracy in the decision trials. Thus, 27 participants  
619 completed session 2. One of these participants was excluded due to a technical error with MEG data collection.  
620 The final sample consisted of 26 right-handed participants (8 males, 18 females) aged between 18 and 35 years  
621 (M = 25, SD = 5).

622 All participants were fluent or native English speakers with normal vision and no current use of psychiatric  
623 medication. Each participant provided written consent for each session and were paid £50 (£10 for behavioural  
624 session and £40 for MEG session), plus up to £15 bonus (up to £5 for the behavioural session and up to £10 for  
625 the MEG session) upon completing the study. Bonuses were calculated by converting the accuracy of each block  
626 (i.e., the proportion of times participants made the correct choice) into a monetary value between £0 and £1.

## 627 Experimental task

### 628 *Image learning*

629 The experiment was created for web browser using jsPsych v6.1.0. The experiment was presented in the format  
630 of a computer game where participants played the role of an astronaut exploring rooms within a spaceship. There  
631 were six rooms in total, arranged as two sequences (or “paths”): path 1 contained rooms A, B, and C, and path 2  
632 contained rooms D, E, and F. Each room (or “state”) was represented by a unique image randomly selected from  
633 a set of 12 for each participant (**Extended Data Fig. 1A**). During the image learning phase, participants watched  
634 an animation of the transitions along each path, in which the images for each room were presented one at a time  
635 for 3 seconds each (**Extended Data Fig. 1C**). Participants were then tested on their memory for the order of images  
636 in each sequence. Participants were given up to two attempts to reach at least 80% accuracy.

### 637 *Value learning*

638 After successfully completing the image learning phase, participants then learned to associate an integer value  
639 (ranging from -5 to 5, excluding 0) with each room. This integer represented the number of points subjects stood  
640 to gain or lose in each room. To learn these values, participants were presented with each sequence four times,  
641 with the integer value presented underneath each image (4-second presentation; **Extended Data Fig. 1D**).  
642 Participants were then tested on their memory for each individual room’s value, as well as their ability to calculate  
643 the cumulative sum of points in each room. This process was repeated until participants scored at least 80%  
644 accuracy (up to two attempts).

645 *Decision trials*

646 After completing image and value learning, participants then partook in decision trials. At the beginning of each  
647 decision trial, participants were placed conceptually “outside” of the environment containing the two learned  
648 sequences and could choose to either approach or avoid it (**Fig. 1A**). Avoidance resulted in a guaranteed point  
649 increase of +1 and no transition to either path. Approach decisions took participants down one of the two paths,  
650 as chosen by the computer. Crucially, however, there was always a degree of uncertainty as to which of the two  
651 paths the participant would transition to if an approach decision was made. The transition probability of each  
652 path varied from trial to trial and was explicitly conveyed to the participant at the beginning of each trial. There  
653 were five possible sets of probabilities: 10-90%, 30-70%, 50-50%, 70-30%, and 90-10% for transitions to paths 1  
654 and 2, respectively. Once transitioned to a path, the transitions to each room within the sequence were  
655 deterministic.

656 Participants were required to use the value map they had learned in the previous stage, in conjunction with the  
657 path transition probabilities presented on each trial, to evaluate the utility of making an approach versus an avoid  
658 decision. Optimally, this evaluation would reflect an expected value calculation for both approach and avoid  
659 choices, such that:

$$EV_{app} = P_1R_1 + P_2R_2 \quad (1)$$

660 where  $EV_{app}$  is the expected value of approaching,  $P_1$  and  $P_2$  are the probabilities of transitioning to paths 1 and 2,  
661 respectively, and  $R_1$  and  $R_2$  are the total sums of points for paths 1 and 2, respectively, taking into account the  
662 odd rule states (see Methods, Planning manipulation). The expected value of avoiding,  $EV_{av}$ , was always 1. The  
663 decision to approach was considered correct if  $EV_{app} \geq EV_{av}$  and the decision to avoid was considered correct if  
664  $EV_{app} \leq EV_{av}$ .

665 After each block, the proportion of correct responses was converted into a monetary value and displayed as a  
666 bonus. Participants did not receive feedback on the accuracy of their choices throughout the block. They did,  
667 however, observe an animation of their subsequent transitions and change in points (**Fig. 1A**). For “avoid”  
668 decisions, a screen was displayed with text stating that they had received 1 point (3 seconds). For “approach”  
669 decisions, participants were first shown which path had been selected by the computer according to the transition  
670 probability (“Path 1” or “Path 2”, for 3 seconds). Participants were then shown each state within that path one at  
671 a time (2 second presentation), underneath which the state value as well as the running total of points collected  
672 along the path was displayed. A blank screen was presented between states (randomly jittered duration between

673 0.5 and 0.8 seconds). A final screen conveyed the total number of points earned for that trial (2 seconds). Trials  
674 were separated by a blank screen (1 second).

675 After an initial practice block, participants completed 6 (behavioural session) or 10 (MEG session) blocks. Each  
676 block contained 18 decision trials. In the practice block, participants were given unlimited time to make their  
677 decision and did not earn bonus money. In test blocks, participants were given 30 seconds (indicated by an on-  
678 screen timer) to make their choice. Responses were disabled for the first 5 seconds to prevent accidental presses  
679 and encourage planning. If no response was made after 30 seconds, participants were penalised -1 point and  
680 prompted with a warning message ("Too slow!") and the trial ended.

### 681 *Planning manipulation*

682 A number of additional features were incorporated into the design of the decision trials to encourage planning, as  
683 well as to control for certain variables. One feature was what we term the "odd rule". The purpose of the odd rule  
684 was to allow the sum of points along each path to vary from trial to trial, thus encouraging participants to engage  
685 in sequential planning. On each trial, the odd rule was applied to two states: one from each path. These two odd  
686 rule states were displayed on-screen (as images on forced-choice trials or as text labels on free-choice trials) at  
687 the beginning of each trial, alongside the path probabilities (**Fig. 1A**). Participants were instructed that, if the sum  
688 of points accumulated up until (and including) an odd rule state was an odd number, then the sign of this  
689 cumulative sum would "flip" (i.e., a negative cumulative sum will become positive, and vice versa). This new sum  
690 would then be carried over to any subsequent states along the path.

691 By way of example, assume the values of states A, B, and C in path 1 are -5, -2, and 3, respectively. If state B is the  
692 odd rule state, then one must mentally sum the number of points up until (and including) state B ( $-5 + -2 = -7$ ). One  
693 must then consider whether the current sum of points is an odd number. In this case, it is (-7), and thus the sign  
694 of the sum is flipped (becoming +7). This value is then carried over to the next state, C ( $7 + 3 = 10$ ), producing a  
695 final outcome of 10. If, instead, state C is the odd rule state, then one sums the number of points up until state C  
696 ( $-5 + -2 + 3 = -4$ ). In this case, the sum of points at the odd-rule state is an even number (-4), and thus no sign-  
697 flipping occurs, producing a final outcome of -4. Hence, the final value of each path is entirely dependent on the  
698 position of the odd rule state in each path (see **Fig. 1A** for another example). This manipulation increased the  
699 variability of final path values across trials. In the MEG session, participants were instructed to refrain from  
700 verbalising numbers aloud to minimise movement-related artefacts in the MEG activity.

701 To further increase the variability in final path values across the experiment, the value of one state from each path  
702 changed at the beginning of each block. All state values then remained constant for the duration of the block. So  
703 that participants knew which values had changed at the beginning of each block, the first four trials in each block  
704 (first six in the practice block) were forced-choice, such that participants could only choose to approach. Forced-

705 choice trials were controlled so that they lead to an equal number of transitions to path 1 and path 2. Any points  
706 gained or lost on these trials did not count towards bonus payment and were not included in planning-related  
707 MEG analyses, as participants were unable to plan until having observed the updated values in both paths.

708 Forced-choice trials were also the only trials in which the images were displayed, both during the planning period  
709 (where the states with the odd rules were displayed) and the sequence animation. In all other free-choice trials,  
710 images were replaced by their text labels (e.g., “cat” or “bicycle”), which had already been shown to participants  
711 during the functional localiser (see **Procedure** below). This was done to control for any potential biased visual  
712 exposure to the state images during free-choice trials based on choice behaviour (e.g., only deciding to approach  
713 when path 1 is more likely) while still periodically reminding participants of the images associated with each room.

714 Participants were assigned to one of two experimental protocols in a counterbalanced fashion (**Extended Data**  
715 **Fig. 1E**). Each protocol was designed to minimise the repetition of odd rule state pairs across trials. These two  
716 protocols also captured another feature of the design, in which one path more often resulted in a positive outcome  
717 and the other in a negative outcome. This was done to maximise the difference in replay between rewarding and  
718 aversive paths, by allowing for some degree of association by repetition. To prevent participants from relying on  
719 this consistency (and thus not engaging in sequential planning), 5% of trials were catch trials, where either both  
720 paths produced a gain or both produced a loss, thus increasing the utility of planning on every trial. Furthermore,  
721 the rewarding and aversive paths swapped positions halfway through the experiment (e.g., if path 1 was  
722 consistently rewarding at the beginning, it became consistently aversive, and vice versa for path 2). The starting  
723 positions of the rewarding and aversive paths were counterbalanced across the two protocols.

## 724 Procedure

### 725 *Initial session*

726 Participants completed two sessions on consecutive days. The first session was a behavioural-only practice,  
727 where participants completed three questionnaires: the 12-item Intolerance of Uncertainty Scale<sup>63</sup>, 16-item Penn  
728 State Worry Questionnaire<sup>64</sup>, and 30-item Domain-Specific Risk-Taking Scale<sup>65</sup>, each presented in a random order  
729 on a computer (approximately 15 minutes). Participants then completed a shorter 45-minute version of the  
730 experiment. The aim of this session was to ensure participants were capable of performing the task (at least 80%  
731 performance on the image and value memory tests, and at least 60% correct choices on decision trials) before  
732 continuing to the MEG session the following day.

## 733 *Functional localiser*

734 The second session comprised an MEG session. Participants first completed a functional localiser task (30-  
735 minutes) and then completed a full 1.5-hour task. In the functional localiser, participants were shown the six  
736 unique images (randomly selected per participant) used in the main task. Crucially, these images were different  
737 from those shown in the initial behavioural session. On each trial, an image was presented on screen for 1 second  
738 (**Extended Data Fig. 1B**). After the image disappeared, two words were presented on the left and the right of the  
739 screen. One of these was the correct label for the previous image (e.g. "cat") and the other label was randomly  
740 selected from a pool of invalid words. Participants pressed either the left or right button of a 4-button response  
741 pad to indicate the correct label. After making a response, the words were replaced by a fixation cross for a  
742 randomly jittered inter-trial interval between 0.5 and 1.5 seconds. Correct and incorrect responses produced a  
743 green or red cross, respectively. There were four blocks, within which each image was randomly presented 20  
744 times, giving 80 trials in total per image. Across the 26 participants, the mean response accuracy was 97.48% (SD  
745 = 2.48%, range = 90.63 to 99.79%).

## 746 MEG analysis

### 747 *MEG acquisition and preprocessing*

748 Participants' neural activity was measured using a CTF Omega MEG scanner with a 275-channel axial gradiometer  
749 whole-head system (CTF Omega, VSM MedTech) at University College London. Participants were seated upright  
750 in the scanner and head position was continuously monitored by three head position indicator coils located at the  
751 nasion and left and right pre-auricular fiducial points. Data were acquired continuously at 1,200 Hz and  
752 participants' eye movements were recorded using an Eyelink eye-tracking system. Triggers were recorded using  
753 a photodiode positioned behind the stimulus presentation screen that detected the onset of a flashing white  
754 stimulus (hidden from view) that was synchronised with event onsets.

755 MEG data from the functional localiser and decision trials were preprocessed using SPM12 (Wellcome Centre for  
756 Human Neuroimaging), Fieldtrip (2019), and custom code written in MATLAB R2018b (MathWorks). All code is  
757 available on GitHub: <https://github.com/jjmcfadyen/approach-avoid-replay>. CTF data for each block were  
758 imported using OSL (the OHBA Software Library, from OHBA Analysis Group). Trigger onset times and durations  
759 were extracted from the photodiode signal and semi-automatically checked for errors. Next, the data were high-  
760 pass filtered at 0.5 Hz to reduce slow drift, and a notch filter for 50 Hz was applied to remove line frequency. The  
761 data were then downsampled to either 100 Hz (for replay analysis, to reduce temporal autocorrelation) or 600 Hz  
762 (for source reconstruction), thereby reducing computational load and increasing signal to noise ratio. OSL also



763 identified potential bad channels whose characteristics fell outside the normal distribution of values for all  
764 sensors.

765 Independent component analysis was then performed on the data (FastICA,  
766 <http://research.ics.aalto.fi/ica/fastica>), decomposing it into 150 independent spatiotemporal components.  
767 Artefactual components were automatically classified using the combined spatial topography, time course, time  
768 course kurtosis, and frequency spectrum of all components. For example, eye blink artifacts exhibited high  
769 kurtosis (>20), a repeated pattern in the time course, and consistent spatial topographies. The number of excluded  
770 components was limited to a maximum of 20. Artefacts were rejected by subtracting them out of the data. All  
771 subsequent analyses were performed directly on the filtered, cleaned MEG signal, in units of femtotesla.

772 The data were then divided into different epochs using the trigger onsets and durations. For the functional  
773 localiser, epochs were created for the image onset (-0.1 to 0.8 seconds post-stimulus onset). For the decision  
774 trials in the main task, epochs were created for the planning time (-0.1 seconds before trial onset to the response  
775 time). Artefactual sensors identified by OSL were interpolated for all epochs, and artefactual functional localiser  
776 trials were excluded from the classification procedure.

### 777 *Image classification*

778 We used Temporal Delayed Linear Modelling (TDLM) to characterise patterns of neural dynamics during the  
779 task<sup>66</sup>, as performed in previous studies<sup>12,14,15,17,47</sup>. First, for each participant, we classified patterns of multivariate  
780 neural activity evoked by each image in the functional localiser (**Fig. 2A**). The purpose of these classifiers was to  
781 detect reinstatement of each image representation during planning, likely indicating memory reactivation. This  
782 approach capitalises on the similarity between spatial patterns of neural activity evoked by the visual onset of  
783 stimuli during conscious viewing and memory retrieval, which has previously been demonstrated in both MEG and  
784 fMRI<sup>67,68</sup>. Specifically, an interplay between hippocampus and distributed cortical networks during memory  
785 retrieval produces spatial patterns of activity that closely resemble patterns of activity that were produced when  
786 stimuli were first experienced<sup>69</sup>. Notably, our stimuli were visually and categorically unique, thus maximising our  
787 ability to detect features reinstated during planning (e.g., visual imagery, conceptual associations, etc.)<sup>47</sup>.

788 We selected data from 0 to 300 ms from each functional localiser epoch, excluding incorrect and artefactual trials,  
789 as well as trials where response time was > 5 standard deviations from the mean per participant (average of 78  
790 trials per stimulus, per participant; SD = 2, range = 73 to 80). We then constructed a series of Lasso-regularised  
791 logistic regression models. Each model received data from a single time sample (0 to 300 ms, at 10-ms resolution)  
792 across all trials. Hence, we constructed separate models (per time sample, and per image; 31 × 6) per participant,  
793 each using a trials × sensors (e.g., 480 × 275) data matrix and a binary vector indicating which trials belonged to

794 that image. For each model, we appended a duplicate-sized matrix of zeros to the data matrix to reduce the spatial  
795 correlation between each model.

796 Each lasso-regularised logistic regression model used a range of 100 regularisation parameters ( $\lambda$ ) sampled from  
797 a half-Cauchy distribution ( $\gamma = 0.05$ , range = 0.0001 to 1). Thus, each model produced a  $\lambda \times$  sensors (100  $\times$  up to  
798 275) matrix of slope coefficients (**Fig. 2B**), as well as a vector of intercept coefficients for each  $\lambda$ . We refer to  
799 these coefficients as our binomial classifiers, each of which are trained to distinguish the sensor data associated  
800 with one image as compared to all other images.

801 To evaluate the accuracy of each classifier per participant, we conducted a  $K$ -folds cross-validation procedure.  $K$   
802 was set to the minimum number of trials per stimulus for that participant. In each fold, a test set was created by  
803 randomly taking one sample from one exemplar trial per stimulus. The remaining data was used for training.  
804 Random selection of the test data was controlled to maximise equal sampling across trials. The classifiers per  
805 state generated from the training dataset were then applied to the six test trials (one for each stimulus). Thus, for  
806 a given fold, a score of 1 or 0 was given for whether each state classifier maximally predicted the correct trial.  
807 The accuracy of each state classifier was given by the average score across folds.

808 For each subject, we selected  $\lambda$  that produced the highest mean accuracy across state classifiers ( $\lambda$ :  $M = 0.0017$ ,  
809  $SD = 0.0015$ ). We then averaged the classification accuracy across states per subject and examined which  
810 training times produced the highest accuracy across subjects (**Fig. 2C**). Overall average state classification  
811 accuracy exceeded chance (16.66%) for all subjects from 80 ms onwards, peaking at 120 ms (48.97%). Classifier  
812 training times from 110 to 150 ms made up the top 15% performance (all > 45.80% accuracy).

### 813 *Sequential state reactivation*

814 Using our state classifiers, we then estimated the degree to which images were sequentially reactivated in the  
815 brain while participants planned whether to approach or avoid the state space in each trial. We utilised an updated  
816 general linear modelling approach, which encapsulates a lagged cross-correlation between the evidence for state-  
817 to-state transitions. This method produces an overall “sequenceness” statistic at different time intervals, or “lags”.  
818 We employed this approach on a trial-by-trial basis per participant, using neural data collected during the planning  
819 period.

820 In the first step, we estimated the degree to which each state was reactivated during the planning period of free-  
821 choice decision trials by multiplying the spatiotemporal MEG data by each state classifier’s beta estimates. We  
822 used state classifiers trained at 120 ms post-stimulus onset, which had the highest cross-validated accuracy  
823 across subjects. We then entered the resultant time series of predicted state reactivation (states  $\times$  time matrix;

824 Fig. 2D) per trial into a 2-level general linear model designed to test whether reactivation of each stimulus occurred  
 825 in a specific order at different time intervals.

826 At the first level, we performed a family of multiple regressions for each state's reactivation time series ( $i \in [1: 6]$ ),  
 827 in which a time-lagged copy of the reactivation time series for state  $j$  ( $X(t\Delta)_j$ ) predicts the original, unshifted  
 828 reactivation time series of state  $i$  ( $X_i$ ). The time lags ranged from 0 to 600 ms, in 10 ms bins. Hence, this analysis  
 829 evaluated the average likelihood that stimulus  $i$  is followed by stimulus  $j$  after a time lag of  $t\Delta$ . Separate linear  
 830 models were estimated for each stimulus  $i$  and each time lag  $t\Delta$ :

$$X_i = \sum_{j=1}^6 X(t\Delta)_j \times \beta(t\Delta)_{ij} + C \quad (2)$$

831 where  $C$  is a constant term and  $\beta(t\Delta)_{ij}$  is a coefficient derived from ordinary least-squares that captures the  
 832 unique influence of  $X_j$  on  $X_i$ . These coefficients are then used to form  $6 \times 6$  empirical transition matrices,  
 833  $\beta(t\Delta)$ , for each time lag.

834 At the second level, we quantified the evidence for specific, hypothesised state-to-state transitions. In this task,  
 835 the key state-to-state transitions were  $A \rightarrow B$  and  $B \rightarrow C$  (path 1), as well as  $D \rightarrow E$  and  $E \rightarrow F$  (path 2). These  
 836 transitions were declared by separate  $6 \times 6$  binary matrices for hypothesised forward ( $T_F$ ) and backward ( $T_B$ )  
 837 transitions, where  $T_F = T_B'$ . The evidence for the hypothesised transitions was then quantified by:

$$B(\Delta t) = \sum_r Z(r) \times T_r \quad (3)$$

838 where  $r$  is the total number of all regressors included in the second level. These regressors included  $T_F$ ,  $T_B$ ,  
 839  $T_{auto}$  (an identity matrix of self-transitions to control for autocorrelation), and  $T_{const}$  (a constant matrix that  
 840 models away the average of all transitions, ensuring that any weight on  $T_F$  and  $T_B$  was not due to general dynamics  
 841 in background neural dynamics). Note that there were four versions of  $T_F$  and  $T_B$ , one for each hypothesised  
 842 transition ( $A \rightarrow B$ ,  $B \rightarrow C$ ,  $D \rightarrow E$ , and  $E \rightarrow F$ ). This allowed us to examine the evidence of replay of each transition  
 843 specifically, which was critical to our path-specific analyses.  $Z$  is the weight for each regressor, representing the  
 844 evidence for the hypothesised state-to-state transitions.  $Z_F$  and  $Z_B$  are evidence for forward and backward  
 845 transitions, respectively. A forwards-minus-backwards sequenceness measure,  $Z_D$ , was also computed by  
 846 performing  $Z_F - Z_B$ , thus removing common variance. Repeating equation 3 at each time lag produces a time

847 series of sequenceness at different intervals, where smaller intervals indicate more time-compressed replay (Fig.  
848 2E).

849 To determine the statistical significance of  $Z$  (averaged across the four transitions and all trials per participant),  
850 we employed non-parametric permutation testing at the second level. We generated a null distribution by  
851 generating all possible invalid versions of  $T_F$  and  $T_B$ , such that they only included cross-path transitions (e.g., A to  
852 E, B to D, etc.). This produced 40 null versions of  $Z$ . We then calculated a significance threshold for our valid  $Z$  by  
853 taking the maximum absolute value of each null and computing the 95th percentile for  $Z_F$  and  $Z_B$  (one-sided test)  
854 or the 2.5th and 97.5th percentile for  $Z_D$  (two-sided test). Thus, values of  $Z$  were deemed statistically significant  
855 (FWE < 0.05) if they exceeded these significance thresholds.

856 To account for inter-subject variability in classification accuracy across training times and their relevance to  
857 replay, we also computed sequenceness using classifiers trained on 110 ms and 130 ms (10 ms either side of  
858 the winning training time). Thus, we computed sequenceness three times per subject, and chose the classifier  
859 training time (110 ms, 120 ms, or 130 ms) that produced the greatest absolute value of  $Z_D$  across lags, averaged  
860 across all transitions (110 ms = 11 subjects, 120 ms = 10 subjects, 130 ms = 5 subjects).

## 861 *Source localisation*

862 To investigate the neural sources underlying replay during planning, we used a procedure for identifying replay  
863 onsets similar to previous studies<sup>14,15,47</sup>. Replay onsets were defined as time samples where reactivation of one  
864 state was followed by reactivation of the following state to a higher degree than that expected by chance. For  
865 each trial, we multiplied the state reactivation matrix ( $X$ ) by a time-shifted version of the state reactivation matrix  
866 by lag  $t$  ( $X(\Delta t)$ ). We did this separately for each lag found to be significant in the group-level replay analysis (20  
867 to 90 ms) and only investigated forward transitions, as only forwards replay was significant at the group level  
868 (**Fig. 3E**). Next, we multiplied  $X$  by a state transition matrix ( $P$ ) that either represented the true sequential order of  
869 states or a randomised order (40 randomisations in total, matching the null iterations used in the replay analysis).  
870 Then, for each lag and for each iteration of  $P$ , we multiplied  $X(\Delta t)$  by  $P$  to produce a matrix of sequential state  
871 reactivation (i.e., replay) per transition across time. We then summed across transitions to produce a vector ( $R$ )  
872 reflecting an overall estimate of replay.

873 To demarcate the onset of a replay event, we estimated a significance threshold in a similar manner to the replay  
874 analysis. For each null iteration of  $P$ , we concatenated  $R$  vectors for all lags and all trials into a single vector, which  
875 were combined to create matrix  $N$  (40 columns: one per randomised state order). We then calculated a  
876 significance threshold by computing the maximum value across the columns of matrix  $N$ , and then computing its  
877 95th percentile. Thus, this permutation approach controlled for multiple comparisons across time samples and  
878 lags, and also maximised our ability to distinguish signal from noise. Individual replay events were marked as

879 instances where replay at any lag exceeded the overall significance threshold. Finally, we excluded any replay  
880 events that were preceded by another replay event (of any lag) in the preceding 100 ms.

881 We epoched the MEG data according to the replay onsets (-100 to 150 ms surrounding replay onset) and baseline  
882 corrected the data using a -100 and -50 ms window. We then transformed these data to a three-dimensional grid  
883 in MNI space (grid step = 5 mm) using a linearly constrained minimum variance beamformer<sup>70,71</sup>, as implemented  
884 in OSL. Forward models were generated on the basis of a single shell using superposition of basis functions that  
885 approximately corresponded to the plane tangential to the MEG sensor array. The sensor covariance matrix for  
886 beamforming was estimated using data separately in theta (4 to 8 Hz) and high gamma (120 to 150 Hz) frequency  
887 ranges.

888 At the first level, we computed one-sample tests on whole-brain source activity at each time point using  
889 nonparametric permutation testing<sup>72</sup> as implemented in OSL. We selected the resultant t-maps for each  
890 participant and smoothed the images in SPM12 using a 12 mm FWHM Gaussian kernel. We then entered these  
891 into one-sample t-tests (averaged from 0 to 100 ms post-replay onset) in SPM12 for group-level inference, with  
892 or without participant trait anxiety or overall performance accuracy added as a covariate. All statistics are  $p <$   
893 0.05, FWE-corrected at the whole brain cluster level. Anatomical labelling was determined via the Automated  
894 Anatomical Labelling Atlas (AAL3) add-on to SPM12<sup>73</sup>.

## 895 Multi-level modelling

896 All analyses were conducted on the MEG session, as the initial behavioural session served purely to acquaint  
897 participants with the task structure. We adopted a multi-level modelling approach, which allowed us to examine  
898 effects on a trial-by-trial basis. This approach also allowed us to compare conditions with unbalanced trial  
899 numbers (e.g., “approach” decisions mostly consisted of trials where reward probability was high, and vice versa  
900 for “avoid” decisions).

901 We used the lme4 package implemented in R v3.6. We constructed a series of models that either used: a) choice  
902 as a binomial dependent variable, or b) sequenceness as a linear dependent variable. In all models, forced choice  
903 trials and catch trials (i.e., trials where both paths resulted in an overall loss or both resulted in an overall gain)  
904 were excluded. All predictors were mean-centred. To ensure convergence, the bobyqa optimiser was used and  
905 set to  $10^6$  iterations. Significant interaction terms were followed up by simple slopes analyses using the  
906 “interactions” package in R, FDR-corrected for multiple comparisons, and the “emmeans” package in R. We also  
907 ensured that all models produced a variable inflation factor (VIF) below 5 and that autocorrelation within the  
908 residuals of each model was minimal, as assessed by a Durbin-Watson test<sup>74</sup>; see **Supplementary methods**).

909 For models including individual differences, we used principal components analysis to reduce the dimensionality  
910 of the three self-report questionnaires (intolerance of uncertainty, worry, and risk-taking across 7 domains: ethical,  
911 social, health, financial, and recreational) completed at the beginning of the behavioural session. We identified  
912 two principal components that together explained 60.55% of the variance (41.52% and 29.49%, respectively;  
913 eigenvalues = 1.548, 1.357, 1.040, 0.944, 0.684, 0.457, 0.350). The first component mapped positively on to risk-  
914 taking questionnaire scores, while the second component mapped negatively on to intolerance of uncertainty and  
915 worry. We refer to these two components as risk-seeking and anxiety, respectively. For interpretability, we inverted  
916 these factors, such that more positive values represented higher risk-aversion and higher anxiety, respectively.

## 917 Statistics and reproducibility

918 No statistical methods were used to pre-determine sample sizes but our sample size was similar to those reported  
919 in previous publications observing significant replay of state transitions as measured with MEG <sup>12,15</sup>, as well as a  
920 relationship between replay and individual differences in performance <sup>47</sup>. All statistical analyses were performed  
921 using computer code available online (see **Code availability**). Raw behavioural and MEG data are also available  
922 online in the interest of experimental reproducibility (see **Data availability**). As the study was a within-subjects  
923 design, there was no randomisation to experimental conditions and thus no blinding during data collection or  
924 analysis. The stimuli presented to each participant was, however, randomised using a random seed generator  
925 based on computer time at the beginning of the experiment. Assumptions of all tests were formally tested. In  
926 cases where assumptions of normality were violated, data were log-transformed. Two participants' data were  
927 excluded from path-specific MEG replay analysis due to poor behavioural performance in the task (< 60%  
928 accuracy), meaning that these participants were unlikely to have processed that the two paths in the experiment  
929 resulted in an overall reward or loss.

## 930 Methods-only references

- 931 63. Carleton, R. N., Norton, M. A. P. J. & Asmundson, G. J. G. Fearing the unknown: a short version of the  
932 Intolerance of Uncertainty Scale. *J. Anxiety Disord.* **21**, 105–117 (2007).
- 933 64. Meyer, T. J., Miller, M. L., Metzger, R. L. & Borkovec, T. D. Development and validation of the Penn State  
934 Worry Questionnaire. *Behav. Res. Ther.* **28**, 487–495 (1990).
- 935 65. Blais, A.-R. & Weber, E. U. A Domain-Specific Risk-Taking (DOSPERT) Scale for Adult Populations. (2006).
- 936 66. Liu, Y. *et al.* Temporally delayed linear modelling (TDLM) measures replay in both animals and humans.  
937 *eLife* vol. 10 Preprint at <https://doi.org/10.7554/elife.66917> (2021).
- 938 67. Kurth-Nelson, Z., Barnes, G., Sejdinovic, D., Dolan, R. & Dayan, P. Temporal structure in associative retrieval.  
939 *Elife* **4**, (2015).
- 940 68. Polyn, S. M., Natu, V. S., Cohen, J. D. & Norman, K. A. Category-specific cortical activity precedes retrieval

- 941 during memory search. *Science* **310**, 1963–1966 (2005).
- 942 69. Horner, A. J., Bisby, J. A., Bush, D., Lin, W.-J. & Burgess, N. Evidence for holistic episodic recollection via  
943 hippocampal pattern completion. *Nat. Commun.* **6**, 7462 (2015).
- 944 70. Woolrich, M., Hunt, L., Groves, A. & Barnes, G. MEG beamforming using Bayesian PCA for adaptive data  
945 covariance matrix regularization. *Neuroimage* **57**, 1466–1479 (2011).
- 946 71. Van Veen, B. D., van Dronkelen, W., Yuchtman, M. & Suzuki, A. Localization of brain electrical activity via  
947 linearly constrained minimum variance spatial filtering. *IEEE Trans. Biomed. Eng.* **44**, 867–880 (1997).
- 948 72. Hunt, L. T. *et al.* Mechanisms underlying cortical activity during value-guided choice. *Nat. Neurosci.* **15**, 470–  
949 6, S1–3 (2012).
- 950 73. Rolls, E. T., Huang, C.-C., Lin, C.-P., Feng, J. & Joliot, M. Automated anatomical labelling atlas 3. *Neuroimage*  
951 **206**, 116189 (2020).
- 952 74. Fox, J. *Applied Regression Analysis and Generalized Linear Models*. (SAGE Publications, 2015).