

## The roles of online and offline replay in planning

Eran Eldar<sup>1,2,3,†</sup>, Gaëlle Lièvre<sup>2,3</sup>, Peter Dayan<sup>4,\*</sup>, Raymond J. Dolan<sup>2,3,\*</sup>

1 Departments of Psychology and Cognitive Sciences, Hebrew University of Jerusalem, Israel

2 Max Planck UCL Centre for Computational Psychiatry and Ageing Research, University College London, UK

3 Wellcome Centre for Human Neuroimaging, University College London, UK

4 Max Planck Institute for Biological Cybernetics, Tübingen, Germany

† Correspondence to: [eran.eldar@mail.huji.ac.il](mailto:eran.eldar@mail.huji.ac.il)

\* Equal contribution

### Abstract

Animals and humans replay neural patterns encoding trajectories through their environment, both whilst they solve decision-making tasks and during rest. Both on-task and off-task replay are believed to contribute to flexible decision making, though how their relative contributions differ remains unclear. We investigated this question by using magnetoencephalography to study human subjects while they performed a decision-making task that was designed to reveal the decision algorithms employed. We characterized subjects in terms of how flexibly each adjusted their choices to changes in temporal, spatial and reward structure. The more flexible a subject, the more they replayed trajectories during task performance, and this replay was coupled with re-planning of the encoded trajectories. The less flexible a subject, the more they replayed previously and subsequently preferred trajectories during rest periods between task epochs. The data suggest that online and offline replay both participate in planning but support distinct decision strategies.

## 1 Introduction

2 Online and offline replay are both suggested to contribute to decision making<sup>1-15</sup>, but their  
3 precise contributions remain unclear. Replay of experienced and expected state transitions  
4 during a task, either immediately before choice or following outcome feedback, is  
5 particularly well suited to mediate on-the-fly planning, where choices are evaluated based on  
6 the states to which they lead (this is known as model-based planning). Off-task replay might  
7 serve a complementary role of consolidating a model of a state space, specifying how each  
8 state can be reached from other states and the values of those states. According to this  
9 perspective, both types of replay help subjects make choices that are flexibly adapted to  
10 current circumstances.

11 However, a different possibility is that off-task replay also directly participates in planning,  
12 by calculating and storing a (so-called model-free) decision policy that specifies in advance  
13 what to do in each state<sup>16-19</sup>. Such a pre-formulated policy is inherently less flexible than a  
14 policy that is constructed on the fly, but at the same time it decreases a need for subsequent  
15 online planning when time itself might be limited. Thus, rather than online and offline replay  
16 both supporting the same form of planning, this latter perspective suggests a trade-off  
17 between them. In other words online replay promotes an on-the-fly model-based flexibility,  
18 whereas offline replay establishes a stable model-free policy.

19 Despite the wide-ranging behavioural implications of the distinction between model-based  
20 and model-free planning<sup>20-23</sup>, and much theorising on the role of replay in one or the other  
21 form of planning, to date there is little data to suggest whether online and offline replay have  
22 complementary or contrasting impacts in this regard. Therefore, we tested the relationship  
23 between both online and offline replay and key aspects of decision flexibility that dissociate  
24 model-free (MF) and model-based (MB) planning<sup>24</sup>. For this purpose, we first recorded MEG  
25 signals from human subjects during rest and while they navigated a specially designed state  
26 space. We then characterized each individual subject's

27 flexibility and decision-making algorithm based on task behaviour, and we analysed their  
28 MEG signals seeking evidence of on-task<sup>25</sup> and off-task<sup>10,12</sup> sequences of state  
29 representations.

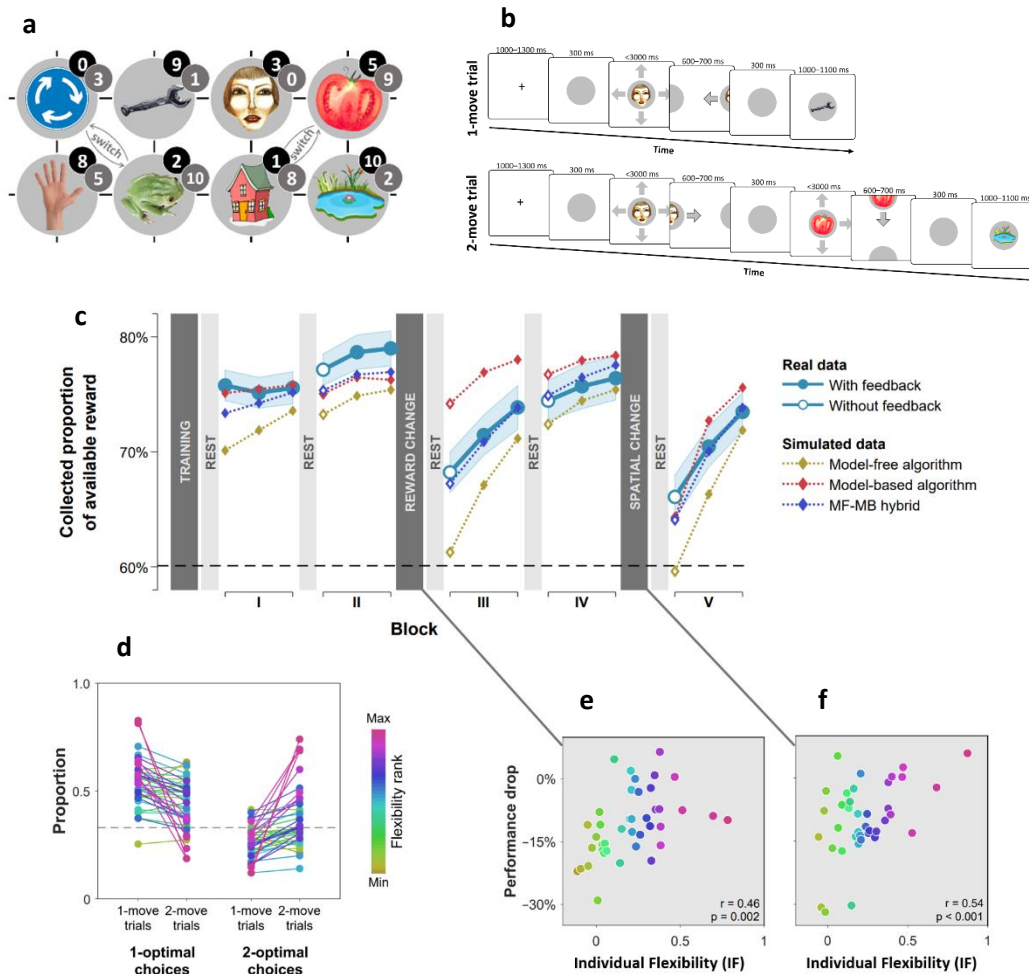
30

## 31 Results

### 32 Individual differences in decision flexibility

33 We used distinct visual images to represent 8 unique states, where occupancy of each state  
34 provided a different amount of reward (**Fig. 1a**). Subjects started each trial at a random state  
35 and had to choose a movement direction in order to collect reward from subsequent states  
36 (**Fig. 1b**). Subjects learnt beforehand how much reward was associated with each state, but  
37 they did not know initially where states were in relation to one another. The latter aspect of  
38 task structure had to be acquired through trial and error learning in order to be able to  
39 implement subsequent moves that delivered the maximal amount of reward.

40 We assessed subjects' flexibility in three ways. First, after the initial two blocks of trials, we  
41 changed the reward associated with each state (**Fig. 1a**; grey numbers) such that persisting  
42 with optimal previous moves would result in below-chance performance. Second, after two  
43 additional blocks of trials, we informed subjects that two specified pairs of states had  
44 switched positions (**Fig. 1a**; 'switch'), again rendering the optimal previous policy now  
45 suboptimal. A flexible model-based planner would be capable of re-planning its moves  
46 perfectly following each of these instructed changes, since such a planner has acquired



**Fig. 1. Subjects differed in decision flexibility.** (a) Experimental task space. Before performing the main task, subjects learned state-reward associations (numbers in black circles) and they were then gradually introduced to the state space in a training session. After performing the main task for two blocks of trials, subjects learned new state-reward associations (numbers in dark gray circles) and then returned to the main task. Before a final block of trials, subjects were informed of a structural task change such that ‘House’ switched position with ‘Tomato’, and ‘Traffic sign’ switched position with ‘Frog’. The bird’s eye view shown in the figure was never seen by subjects. They only saw where they started from on each trial and, after completing a move, the state to which their move led. The map was connected as a torus (e.g., starting from ‘Tomato’, moving right led to ‘Traffic sign’, and moving up or down from the tomato led to ‘Pond’). (b) Each trial started from a pseudorandom location from whence subjects were allowed either one (‘1-move trial’) or two (‘2-move trial’) consecutive moves (signalled at the start of each set of six trials), before continuing to the next trial. Outcomes were presented as images alone, and the associated reward points were not shown. A key design feature of the map was that in 5 out of 6 trials the optimal (first) move was different depending on whether the trial allowed one or two moves. For instance, given the initial image-reward associations (black) and image positions, the best *single* move from ‘Face’ is LEFT (9 points), but when two moves are allowed it is best to move RIGHT and then DOWN (5+9 giving 15 total points). Note that the optimal moves differed also given the second set of image-reward associations. On ‘no-feedback’ trials (which started all but the first block), outcome images were also not shown (i.e., in the depicted trials, the ‘Wrench’, ‘Tomato’ and ‘Pond’ would appear as empty circles). (c) The proportion of obtainable reward points collected by the experimental subjects, and by three simulated learning algorithms. Each data point corresponds to 18 trials (six 1-move and twelve 2-move trials), with 54 trials per block. The images to which subjects moved were not shown to subjects for the first 12 trials of Blocks II to V (the corresponding ‘Without feedback’ data points also include data from 6 initial trials with feedback wherein starting locations had not yet repeated, and thus, subjects’ choices still reflected little new information). All algorithms were allowed to forget information so as to account for post-change performance drops as best fitted subjects’ choices (see **Materials and methods** for details). Black dashed line: chance performance. Shaded area: SEM. (d) Proportion of first choices that would have allowed collecting maximal reward where one (‘1-optimal’) or two (‘2-optimal’) consecutive moves were allowed. Choices are shown separately for what were in actuality 1-move and 2-move trials. Subjects are colour coded from lowest (gold) to highest (red) degree of flexibility in adjusting to one vs. two moves (see text). Dashed line: chance performance (33%, since up and down choices always lead to the same outcome). (e,f) Decrease in collected reward following a reward-contingency (e) and spatial (f) change, as a function of the index of flexibility (IF) computed from panel d. Measures are corrected for the impact of pre-change performance level using linear regression. *p* value derived using a permutation test.

1 knowledge as to how each state can be reached. Conversely, a pure model-free planner would  
2 require complete relearning by trial and error each time there is a change so as to establish a  
3 new policy, since such an agent only possesses a now counterproductive policy that specifies  
4 where to move from each state.

5 Examining how subjects' overall performance altered immediately following these changes  
6 revealed a decrement in average performance (**Fig. 1c**). However, there were substantial  
7 individual differences in this regard, with some subjects seamlessly adapting to reward and  
8 position changes, and others showing drops in performance to chance-levels. Subjects whose  
9 performance showed a strong decline following a reward change tended to cope poorly also  
10 with the position change ( $\rho = 0.50$ , partial correlation controlling for performance levels  
11 before the changes;  $p = 0.001$ , Permutation test).

12 As a third, more continuous, test of a different aspect of decision flexibility, we interleaved  
13 sets of six trials in which only a single move was allowed ('1-move trials') with trials which  
14 allowed two consecutive moves ('2-move trials'; **Fig. 1b**). In 2-move trials, subjects were  
15 rewarded for both states they visited, and thus, an optimal course of action often required  
16 subjects to move first to an initial low-reward state in order to gain access to a high reward  
17 state with their second move. Thus, we defined an individual index (IF) of decision flexibility  
18 as the difference between the proportion of moves that were optimal given the actual number  
19 of allotted moves and the proportion of moves that would have been optimal given a different  
20 number of allotted moves (i.e., had 1-move trials instead involved two moves and 2-move  
21 trials involved one move). A value of zero implies no net adjustment, while positive values  
22 imply advantageous flexibility.

23 The results indicate subjects adjusted their choices advantageously to the number of allotted  
24 moves (+0.21, SEM 0.05,  $p < 0.001$ , Bootstrap test), though there was evidence again of  
25 substantial individual differences (**Fig. 1d**). Importantly, IF correlated with how well a  
26 subject coped with the reward-contingency (**Fig. 1e**) and position (**Fig. 1f**) changes as well  
27 with how accurately they could sketch maps of the state space at the end of the experiment  
28 ( $r = 0.51$ ,  $p < 0.001$ , Permutation test; **Supplementary Fig. 1**). Moreover, examining a  
29 subset of 2-move trials in which subjects made their second moves without seeing the  
30 consequence of their first moves, indicated subjects with high IF planned two steps into the  
31 future (**Supplementary Note 1**), as would be expected from MB planning.

## 32 **Individual flexibility reflected MF-MB balance**

33 These convergent results suggest that IF reflected deployment of a MB planning strategy. To  
34 test this formally, we compared how well different model-free and model-based decision  
35 algorithms, as well as a combination of both, explained subjects' choices. Importantly, we  
36 enhanced these algorithms to maximize their ability to mimic one another (see **Materials and**  
37 **methods** for details). Thus, for instance, the MF algorithm included separate 1-move and 2-  
38 move policies.

39 We found that a hybrid of MF and MB algorithms outperformed substantially either of them  
40 alone (Bayesian Information Criterion<sup>26</sup>: MF = 40821, MB = 43249, MF-MB hybrid =  
41 39908), suggesting that subjects employed a mix of MF and MB planning strategies.  
42 Simulating task performance using the hybrid algorithm showed it captured adequately  
43 differences that were evident between subjects (correlation between real and simulated IF:  
44  $r = 0.92$ ,  $p < 0.001$ , Permutation test; **Supplementary Fig. 2a**). When we examined each  
45 subject's best-fitting parameter values, to determine which of these covaried with IF, we  
46 found 84% of inter-individual variance was explained by three parameters that control a

1 balance between flexible, model-based, and inflexible, model-free, planning (**Supplementary**  
2 **Fig. 2b**). Importantly, less flexible subjects had comparable learning rates and a higher  
3 model-free inverse temperature parameter (in 2-move trials), indicating that lower flexibility  
4 did not reflect a non-specific impairment, but rather, it was associated with enhanced  
5 deployment of a model-free algorithm. Thus, our index of flexibility specifically reflected the  
6 influence of model-based, as compared to model-free, planning.

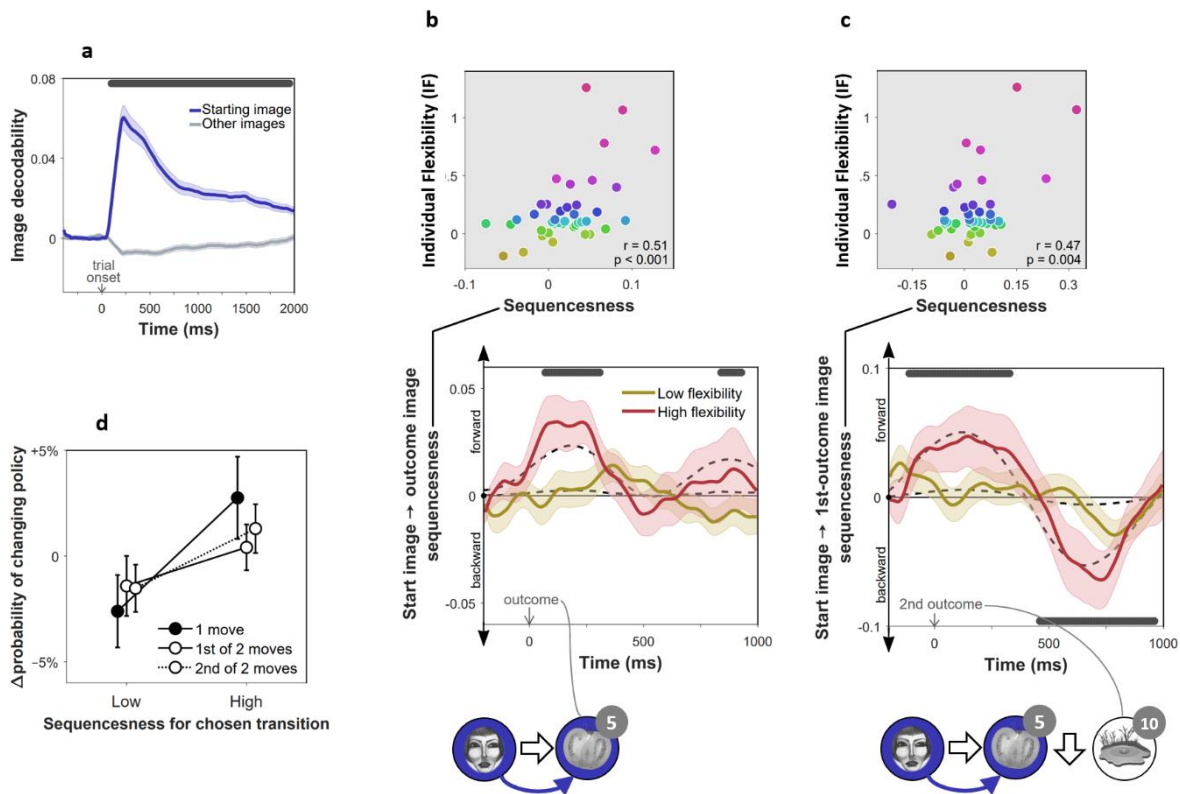
## 7 **On-task replay is induced by prediction errors and associated with high flexibility**

8 In rodents, reinstatement of past states, potentially in the service of planning, is evident both  
9 prior to choices<sup>27</sup> and following observation of outcomes<sup>2</sup>. Thus, we determined firstly at  
10 what point states were neurally reinstated during our task. For this purpose, we trained MEG  
11 decoders to identify the images subjects were processing (**Fig. 2a**). Such decoders robustly  
12 reveal stimulus representations that are reinstated from memory and contribute to decision  
13 processes<sup>25,28</sup>. Crucially, image decoders were trained on MEG data collected prior to  
14 subjects having any knowledge about the task, ensuring that the decoding was free of  
15 confounds related to other task variables (see **Materials and methods**). Applying these  
16 decoders to MEG signals from the main task, we found no evidence of prospective  
17 representation of outcome states (images) to which subjects will transition at choice  
18 (**Supplementary Fig. 4a**). Instead, we found strong evidence that following outcomes  
19 (corresponding to new states to which subjects transitioned), subjects represented the states  
20 from which they had just moved ( $\bar{t} = 3.4$ ,  $p = 0.001$ , Permutation test; **Supplementary Fig.**  
21 **4b**). Consequently, we examined in detail the MEG data recorded following each outcome  
22 for evidence of replay of state sequences that subjects had just traversed.

23 To test for evidence of replay, we applied a measure of “sequenceness” to the decoded MEG  
24 time series, a metric we have previously shown is sensitive in detecting replay of experienced  
25 and decision-related sequences of states<sup>10,12,25</sup>. Importantly, sequenceness is not sensitive to  
26 simultaneous covariation, and thus, it is only found if stimulus representations follow one  
27 another in time<sup>25</sup> (as in previous work, we allowed for inter-stimulus lags of up to 200 ms).  
28 Thus, following each outcome, we computed sequenceness between the decoded  
29 representations of the preceding and the outcome state (**Fig. 2b**). Additionally, MEG signals  
30 recorded following the second outcome in 2-move trials were also tested for sequenceness  
31 reflecting the trial’s first transition (i.e., between the starting state and first outcome; **Fig. 2c**).

32 Using an hierarchical Bayesian Gaussian Process approach (see **Methods** for details) we  
33 tested for timepoints at which sequenceness was evident and correlated with individual  
34 flexibility. This method directly corrects for comparison across multiple timepoints by  
35 accounting for the dependency between them<sup>29</sup>. Since replay is thought to be induced by  
36 surprising observations<sup>16,17,30,31</sup>, we also included surprise about the outcome (i.e. the state  
37 prediction error inferred by the hybrid algorithm) as a predictor of sequenceness. We found  
38 significant sequenceness encoding the last experienced state transition (from 50 to 330 ms  
39 and from 820 to 950 ms following outcome onset; **Fig. 2b**; note that the median split is only  
40 for display purposes; analyses depended on the continuous flexibility index) and, at the  
41 conclusion of 2-move trials, also the penultimate transition (from 130 ms before to 350 ms  
42 following outcome onset; **Fig. 2c**). These sequences were accelerated in time, with an  
43 estimated lag of 130 ms between the images, and were encoded in a ‘forward’ direction  
44 corresponding to the order actually visited. Moreover, later in the post-outcome epoch, the

- 1 penultimate transition was also replayed backwards (from 440 to 940 ms following outcome
- 2 onset).
- 3 Importantly, we found this evidence of replay, across all timepoints, was correlated with IF
- 4 (mean  $\beta = 0.17$ , 95% Credible Interval = 0.13 to 0.20), with surprise about the outcome
- 5 (mean  $\beta = 0.06$ , CI = 0.03 to 0.10) and with the interaction of these two factors (mean  $\beta =$
- 6 0.19, CI = 0.15 to 0.22). Thus, sequenceness was predominantly evident following surprising
- 7 outcomes in subjects with high index of flexibility, consistent with online replay contributing
- 8 to model-based planning.



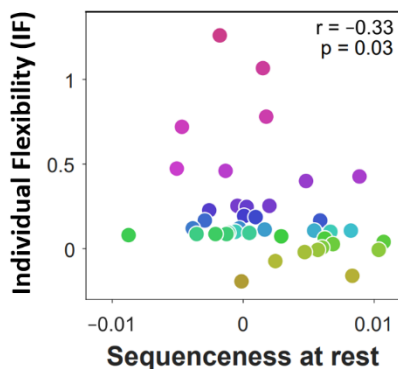
**Fig. 2. On-task replay of state-to-state trajectories as a function of individual flexibility.**  $n = 40$  subjects. (a) Validation of the image MEG decoder used for the sequenceness analyses. The plot shows the decodability of starting images from MEG data recorded during the main task at trial onset. Decodability was computed as the probability assigned to the starting image by an 8-way classifier based on each timepoint's spatial MEG pattern, minus chance probability (0.125). (b) Sequenceness corresponding to a transition from the image the subject had just left ('Start image'; in the cartoon at the bottom, the face) to the image to which they arrived ('outcome image'; the tomato) following highly surprising outcomes (i.e., above-mean state prediction error). In the cartoon, the white arrow indicates the actual action taken on the trial; the blue arrow indicates the sequence that is being decoded. For display purposes only, mean time series are shown separately for subjects with high (above median) and low (below median) IF. Positive sequenceness values indicate forward replay and negative values indicate backward replay. As in previous work<sup>25</sup>, sequenceness was averaged over all inter-image time lags from 10 ms to 200 ms, and each timepoint reflects a moving time window of 600 ms centred at the given time (e.g., the 1 s timepoint reflects MEG data from 0.7 s to 1.3 s following outcome). Dashed lines show mean data generated by a Bayesian Gaussian Process analysis, and the dark gray bars indicate timepoints where the 95% Credible Interval excludes zero and Cohen's  $d > 0.1$ . The top plot shows IF as a function of sequenceness for the timepoint where the average over all subjects was maximal.  $p$  value derived using a permutation test. (c) Sequenceness following the conclusion of 2-move trials corresponding to a transition from the starting image to the first outcome image. (d) Difference in the probability of subsequently choosing a different transition as a function of sequenceness recorded at the transition's conclusion. For display purposes only, sequenceness is divided into high (i.e., above mean) and low (i.e., below mean). A correlation analysis between sequenceness and probability of policy change showed a similar relationship (Spearman correlation:  $M = -0.04$ ,  $SEM = 0.02$ ,  $p = 0.04$ , Bootstrap test). Sequenceness was averaged over the first cluster of significant timepoints from panels b and c, in subjects with non-negligible inferred sequenceness (more than the standard deviation divided by 10;  $n = 25$ ), for the first time the subject chose each trajectory. Probability of changing policy was computed as the frequency of choosing a different move when occupying precisely the same state again. 0 corresponds to the average probability of change (51%).

## 1 On-task replay is associated with changes of policy

2 Recent theorising regarding the role of replay in planning argues that replay should be  
3 preferentially induced when there is a benefit to changing one's policy<sup>17</sup>. This perspective  
4 predicts that, at least in our experiment, subjects should be more disposed to replay  
5 trajectories that they might not want to choose again, rather than trajectories whose choice  
6 reflects a firm policy. To determine whether decodable on-task replay was associated with  
7 policy changes, we tested the relationship between sequenceness corresponding to each move  
8 that subjects chose, and the probability of making a different choice when occupying the  
9 same state later on. We found that moves after which high forward sequenceness was evident  
10 corresponded to moves that were less likely to be re-chosen subsequently (**Fig. 2d**), and these  
11 policy changes increased the proportion of obtained reward ( $M = +11.1\%$ ,  $SEM = 1.5\%$ ,  
12  $p = 0.001$ ). Thus, evidence of online replay was coupled with advantageous re-planning in  
13 relation to the same trajectories.

## 14 Off-task replay is induced by prediction errors and associated with low flexibility

15 We next studied off-task replay, examining MEG data recorded during the 2-minute rest  
16 period that preceded each experimental block. Since each block included five frequently  
17 repeating starting states, we computed sequenceness for the five most frequent image-to-  
18 image transitions subjects chose before and after each rest period (mean choice frequency =  
19 8.4 repetitions per block). As a control analysis, we also examined sequenceness for the five  
20 least frequently chosen transitions from the same starting states (mean choice frequency = 1.0  
21 repetitions per block). We found significant evidence for sequenceness throughout the rest  
22 periods for frequent transitions ( $M = 0.002$ ,  $SEM = 0.001$ ,  $p = 0.01$ , Bootstrap test). By  
23 contrast, no sequenceness was found for the infrequent transitions ( $M < 0.001$ ,  $SEM = 0.001$ ,  
24  $p = 0.47$ , Bootstrap test). Frequent transitions were replayed in a forward direction, with an  
25 estimated time lag of 180 ms between images, and prioritized trajectories that induced more  
26 reward prediction errors in the previous block (correlation of sequenceness with sum of  
27 absolute model-free reward prediction errors inferred by the hybrid algorithm:  $M = 0.04$ ,  
28  $SEM = 0.018$ ,  $p = 0.03$ , Bootstrap test). Most importantly, off-task sequenceness negatively  
29 correlated with IF (**Fig. 3**). This association of sequenceness during rest with low flexibility is  
30 consistent with a proposed role of offline replay in establishing model-free policies<sup>16-19</sup>.



**Fig. 3. Off-task replay of past and future trajectories**.  $n = 40$  subjects. Individual flexibility as a function of sequenceness in rest MEG data for the five most frequently experienced image-to-image transitions. For each rest period, sequenceness was averaged over transitions from both the preceding and following blocks of trials.  $p$  value derived using a permutation test.

## 1 Off-task replay can predict subsequently chosen sequences

2 If offline replay is involved in planning, then its content should predict subjects' subsequent  
3 choices. To test this, we dissociated the replay of experienced trajectories from that of  
4 planned trajectories, focusing on the third rest period after which the optimal image-to-image  
5 transitions changed entirely (due to a change in state-reward associations). As subjects had  
6 been taught about the reward change before this rest period, this afforded an opportunity to  
7 re-plan their choices accordingly during this rest epoch.

8 We first examined the behavioural effect of the state-reward change in more detail. The most  
9 frequently chosen transitions in the block that followed the third rest period differed from the  
10 transitions most frequently chosen in the preceding block (overlap:  $M = 14\%$ ,  $SEM = 3\%$ ),  
11 and this policy change was substantially greater than for the other rest periods (overlap:  $M =$   
12  $53\%$ ,  $SEM = 2\%$ ). As expected, the newly chosen transitions from the following block were  
13 advantageous given the new state-reward associations (reward collected:  $M = 71\%$ ,  $SEM =$   
14  $2\%$ ; chance = 60%) and disadvantageous given the state-reward associations that had so far  
15 applied ( $M = 52\%$ ,  $SEM = 2\%$ ).

16 Given the behavioural change, we focused our examination of the MEG data on evidence for  
17 sequenceness during this crucial third rest period. We found that subjects indeed replayed the  
18 transitions they subsequently chose ( $M = 0.004$ ,  $SEM = 0.002$ ,  $p = 0.02$ , Bootstrap test).  
19 This replay of subsequently chosen moves indicates subjects utilized a model of the task to  
20 re-plan their moves offline<sup>16,17,19</sup>. Our reasoning here is that re-planning in light of the new  
21 reward associations, before subjects experienced them in practice, requires a model that  
22 specifies how to navigate from one state to another. Indeed, multiple regression analysis  
23 showed that low IF was only associated with sequenceness encoding previously chosen  
24 transitions ( $\beta = -0.35$ ,  $t_{37} = 2.25$ ,  $p = 0.03$ ), whereas the replay of subsequently chosen  
25 transitions did not correlate with IF ( $\beta = -0.004$ ,  $t_{37} = 0.03$ ,  $p = 0.97$ ). On the other hand,  
26 the lack of a flexibility enhancement associated with prospective offline replay might indicate  
27 that, as might be expected, offline planning is ill-suited for enhancing trial-to-trial flexibility.

## 28 Discussion

29 We find substantial differences in the behaviour of individual subjects in a simple state-based  
30 sequential decision-making task that correspond also to a distinction in the nature, and  
31 apparent effects, of MEG-recorded on- and off-task replay of state trajectories. These results  
32 bolster important behavioural dissociations, as well as provide substantial new insights into  
33 the control algorithms that subjects employ. The findings fit comfortably with an evolving  
34 literature that addresses human replay and replay<sup>10-12,25,28</sup>.

35 There is an intuitive appeal to the distinction between model-based and model-free reasoning,  
36 confirmed by its close association with many well-established psychological distinctions<sup>32,33</sup>.  
37 However, tasks that have become popular for investigating this distinction<sup>24,34-36</sup> have been  
38 criticized for offering a better grasp on model-based compared to model-free reasoning  
39 processes<sup>36,37</sup>; for rewarding model-based reasoning indifferently<sup>38</sup>; and for admitting  
40 complex model-free strategies that can masquerade as being model-based<sup>39</sup>.

41 In our new task, we show a convergence between superficially divergent methods for  
42 distinguishing model-based and model-free methods – flexibility to immediate task demands  
43 (one-step versus two-step control), preserved performance in the face of changes in the



1 location of rewards or structure, and an ability to reproduce explicitly, after the fact, the  
2 transition structure. Furthermore, the task effectively incentivizes flexible model-based  
3 reasoning, as this type of reasoning alone allows collection of substantial additional reward  
4 (93%) compared to our most successful MF algorithm (80%). These convergent observations  
5 suggest that the model-based and model-free distinction we infer from our task rests on solid  
6 behavioural grounds.

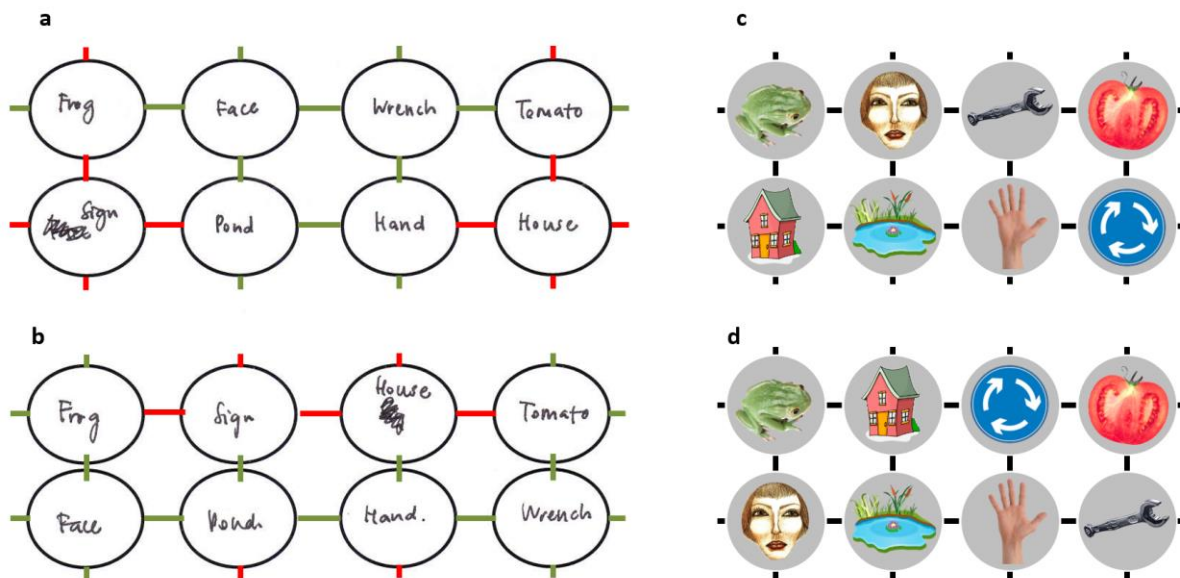
7 In human subjects, there is a growing number of observations of replay and/or preplay of  
8 potential trajectories of states that are associated with the structure of tasks that subjects are  
9 performing<sup>10,11,28</sup>. However, it has been relatively hard to relate these replay events to  
10 ongoing performance. By contrast, there is evidence that rodent preplay has at least some  
11 immediate behavioural function<sup>8,27</sup>, and there are elegant theories for how replay should be  
12 optimally sequenced and structured in the service of planning<sup>17</sup>. In particular, it has been  
13 suggested that replay should prioritize trajectories that can soon be re-encountered, and for  
14 which one's policy can be improved. Our results are broadly consistent with this theoretical  
15 perspective, showing that new surprising observations precede evidence of corresponding  
16 replay, and which in turn predicts appropriate changes in policy. However, rather than  
17 preplay immediately prior to choice, we found evidence of on-task replay following feedback  
18 alone, suggesting a third potential factor impacting on the timing and content of replay – the  
19 need to minimize memory load by embedding new information in ones' policies as soon as it  
20 is received.

21 Critically, the timing and content of replay differed across individuals in a manner that links  
22 with their dominant mode of planning. More model-based subjects tended to replay  
23 trajectories during learning, predominantly reflecting choices they were likely to reconsider.  
24 There have been reports of preferential replay of deprecated trajectories in rodents<sup>8,41</sup>.  
25 However, those studies are consistent with a more general function for replay (e.g.,  
26 maintaining the integrity of a map given a biased experience), whereas in our case, replay  
27 was closely related to future behaviour.

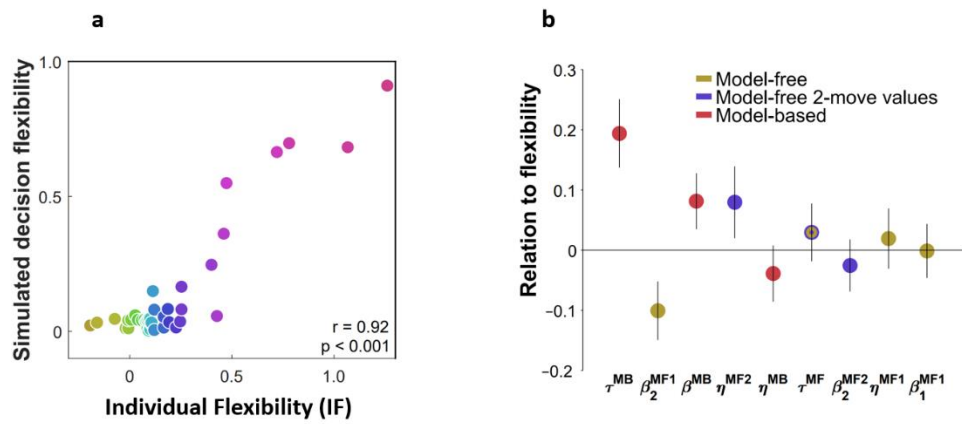
28 By contrast the decodable replay of more model-free subjects centred on rest periods, during  
29 which DYNA-like mechanisms are hypothesized to compile information about the  
30 environment to create an effective model-free policy<sup>17</sup>. This replay of state-to-state transitions  
31 suggests that despite a general inability at the end of the task to draw a map accurately,  
32 model-free subjects do have implicit access to some form of model, though likely an  
33 incomplete one. In any case, generating a policy offline might not be a good strategy for a  
34 task that requires trial-to-trial flexibility, consistent with the lack of association here between  
35 offline replay and ultimate winnings.

36 Our work has a number of limitations. First, our experiment was not ideally suited to  
37 inducing compound representations that link states with those that succeed them, since  
38 succession here frequently changed both within and between blocks. However, algorithms  
39 that utilize such representations mimic both model-free and model-based behaviour, and  
40 future work could utilize our methods to investigate whether and how these algorithms are  
41 aided by online and offline forms of replay<sup>41</sup>. Second, the sequenceness measure that we use  
42 to determine replay suffers from a restriction of comparing forwards to backwards sequences.  
43 There is every reason to expect both forwards and backwards sequences co-exist, so focusing  
44 on a relative predominance of one or the other is likely to provide an incomplete picture. The  
45 problem measuring forwards and backwards replay against an absolute standard is the issue  
46 of a large autocorrelation in the neural decoding, and better ways of correcting for this are  
47 desirable in future studies. Nevertheless, despite these shortcomings the work we report here

- 1 is a further step towards revealing the rich and divergent structure of human choice in
- 2 sequential decision making tasks.



**Supplementary Fig. 1. Example sketches of the state space by a representative subject.** Subjects sketched the state space at the end of the experiment, recalling how it had been structured before and after (b) the position change. On average, subjects sketched much of the state spaces accurately (correct state transitions: first map  $M = 0.65$ ,  $SEM = 0.06$ ; second map  $M = 0.56$ ,  $SEM = 0.06$ ; chance = 0.14,  $p < 0.001$ , Bootstrap test). (a,b) Sketches by a representative subject with 0.58 accuracy for the state space before (a) and after (b) the spatial change. Erroneous transitions are marked in red. (c,d) The actual state spaces the subject navigated before (c) and after (d) the position change.



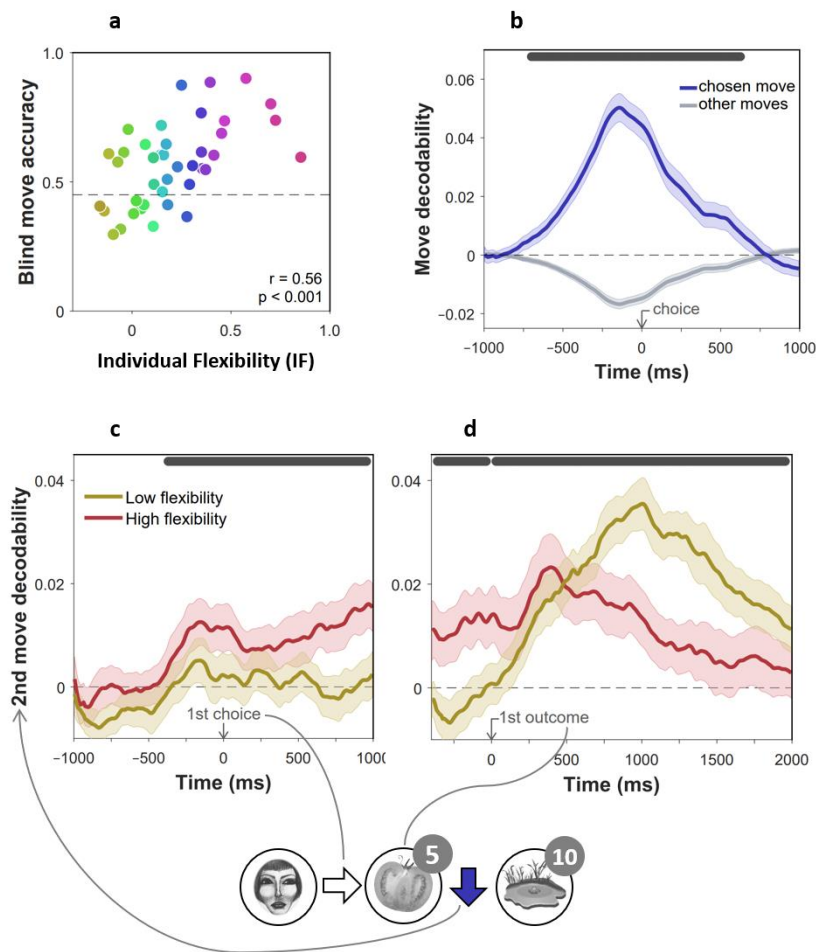
**Supplementary Fig. 2. Individual flexibility reflected the balance between MB and MF planning.**  $n = 40$  subjects. **(a)** Actual and simulated individual flexibility (IF). Task performance was simulated using subjects' best-fitting parameter settings. IF was computed for each simulated subject and averaged over 100 simulations. **(b)** Relationship between IF and individually-fitted parameters. IF was regressed on subjects' best-fitting parameter settings, including all learning ( $\eta$ ), memory ( $\tau$ ), and inverse temperature ( $\beta$ ) parameters. Parameters are color-coded by the component of the algorithm they enhance. Error bars: 95% CI.

## 1 **Supplementary Note 1: Planning two steps into the future**

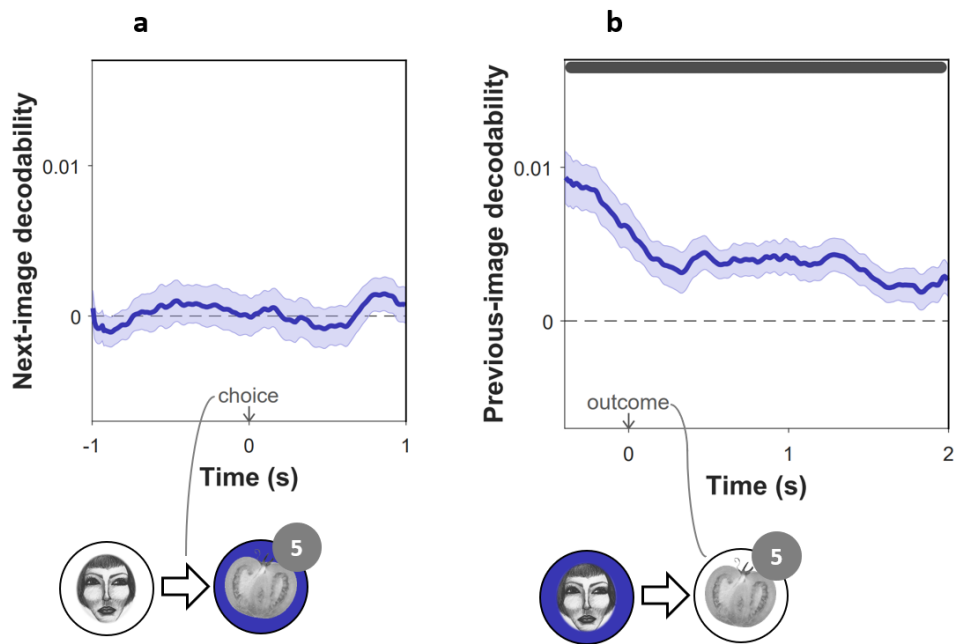
2 Having a cognitive model that specifies how states are spatially organized makes it possible  
3 to plan several steps into the future. To test whether subjects were able to do that, we  
4 challenged subjects with 12 ‘without-feedback’ trials at the beginning of each of the last 4  
5 blocks, during which outcome images were not shown. This meant that in 2-move trials  
6 subjects had to choose their second move ‘blindly’, without having seen the image to which  
7 their previous move had led (e.g., the tomato in **Fig. 1b**). We found that subjects performed  
8 above chance on these blind second moves (proportion of optimal choices: 0.56, SEM 0.03;  
9 chance = 0.45;  $p < 0.001$ , Bootstrap test), and this was the case even immediately following  
10 position and reward changes, when subjects could not have relied on previously tested 2-  
11 move sequences (0.52, SEM 0.03;  $p = 0.01$ , Bootstrap test). Most importantly, such blind-  
12 move success was correlated with IF (**Supplementary Fig. 3a**).

13 This result indicates that more flexible subjects were better able to plan two steps into the  
14 future when required. Examining response times suggested flexibility was associated with  
15 advance planning also when it was not required. Thus, we found that IF correlated with  
16 quicker execution of second moves in general (Spearman correlation with median reaction  
17 time:  $r = -0.61$ ,  $p < 0.001$ , Permutation test). To determine whether advance planning was  
18 indeed generally associated with flexibility, we examined at what point during a trial their  
19 choices became decodeable from MEG signals. For this purpose, we trained a decoder to  
20 decode chosen moves from MEG signals recorded outside of the main task (see **Materials**  
21 **and methods** for details). Validating the decoder on MEG data from the main task showed  
22 that chosen moves became gradually more evident over the course of the trial, their  
23 decodability peaking 140 ms before a choice was made (**Supplementary Fig. 3b**).

24 Thus, we used the move decoder to test whether second-move choices began to materialize in  
25 the MEG signal even before subjects observed the outcomes of their first moves. We found  
26 that chosen second moves were indeed decodeable already during first-move choices  
27 (decodability:  $M = 0.006$ , 95% Credible Interval = 0.004 to 0.008, Bayesian Gaussian  
28 Process analysis; **Supplementary Fig. 3c**) and prior to the appearance of the first outcome  
29 (decodability:  $M = 0.004$ , 95% Credible Interval = 0.002 to 0.006; **Supplementary Fig.**  
30 **3d**). Importantly, this early decodability was correlated with IF ( $\beta$ :  $M = 0.29$ , 95% Credible  
31 Interval = 0.24 to 0.34). By contrast, later decodability, following the onset of the second  
32 image, did not correlate with IF ( $\beta$ :  $M = 0.02$ , 95% Credible Interval =  $-0.02$  to 0.05).  
33 Thus, neural and behavioural evidence concur with the notion that flexibility was associated  
34 with planning second moves prospectively.



**Supplementary Fig. 3. Evidence of advance prospective planning in flexible subjects.**  $n = 40$  subjects. (a) Proportion of optimal choices in second moves for trials without feedback, as a function of individual index of flexibility (IF). In such trials, second moves were enacted without seeing the state they were made from. Measures are corrected using linear regression for accuracy of non-blind moves from the same phases of the experiment. (b) Validation of move decoder. The plot shows the decodability of chosen and unchosen moves from MEG data recorded during the main task. Decodability was computed as the probability assigned to the chosen move (right, left, up or down) by a 4-way classifier based on each timepoint's spatial MEG pattern, minus the average probability assigned to the same moves at baseline (400 ms preceding trial onset). A separate decoder was trained for each subject on MEG data recorded outside of the main task, during the image-reward association training phases. (c,d) Decodability of second moves (the blue arrow in the bottom example cartoon) in 2-move trials during first move choice (c) and presentation of the first outcome (d), as a function of IF. For display purposes only, mean time series are shown separately for subjects with high (above median) and low (below median) IF. In all panels, dark gray bars indicate timepoints where the 95% Credible Interval excludes zero and Cohen's  $d > 0.1$  (Bayesian Gaussian Process analysis). Dashed lines: chance decodability level.



**Supplementary Fig. 4. Previous, not subsequent, images were encoded in MEG.** (a) Decodability during choice, of the image to which the chosen move led subsequently, in high- and low-flexibility subjects. (b) Decodability following outcome, of images subjects had visited earlier in the trial. In both panels, the analysis excluded decoded probabilities assigned to the image presently on the screen. Dark gray bars indicate timepoints where the 95% Credible Interval excludes zero and Cohen's  $d > 0.1$  (Bayesian Gaussian Process analysis). Dashed lines: chance decodability level. Example trials are shown below the plots with decoded elements marked in blue.

## 1 **Materials and Methods**

2 **Subjects.** 40 human subjects, aged 18–33 years, 25 female, were recruited from a subject  
3 pool at University College London. Exclusion criteria included age (younger than 18 or older  
4 than 35), neurological or psychiatric illness, and current psychoactive drug use. To allow  
5 sufficient statistical power for comparisons between subjects, we set the sample size to  
6 roughly double that used in recent magnetoencephalography (MEG) studies on dynamics of  
7 neural representations<sup>28,42</sup>, and in line with our previous study of individual differences using  
8 similar measurements (including ‘sequenceness’)<sup>25</sup>. Subjects received monetary  
9 compensation for their time (£20) in addition to a bonus (between £10 and £20) reflecting  
10 how many reward points subjects earned in the experiment task. The experimental protocol  
11 was approved by the University of College London local research ethics committee, and  
12 informed consent was obtained from all subjects.

13 **Experimental design.** To study flexibility in decision making, we designed a 2x4 state space  
14 where each location was identified by a unique image. Each image was associated with a  
15 known number of reward points, ranging between 0 and 10. Subjects’ goal was to collect as  
16 much reward as possible by moving to images associated with a high numbers of points.  
17 Subjects were never shown the whole structure of the state space, and thus, had to learn by  
18 trial and error which moves lead to higher reward.

19 Subjects were first told explicitly how many reward points were associated with each of the  
20 eight images. Subjects were then trained on these image-reward associations until they  
21 reliably chose the more rewarding image of any presented pair (see **Image-reward training**).

22 Next, the rules of the state-space task were explained (see **State-space task**), and multiple-  
23 choice questions were used to ensure that subjects understood these instructions. To facilitate  
24 learning, subjects were then gradually introduced to the state space, and were allowed one  
25 move at a time from a limited set of starting locations (see **State-space training**). Following  
26 this initial exposure, the rules governing two-move trials were explained and subjects  
27 completed a series of exercises testing their understanding of a distinction between one-move  
28 and two-move trials (see **State-space exercise**). Once these exercises were successfully  
29 completed, subjects played two full blocks of trials in the state space, that included both one-  
30 move and two-move trials.

31 We next tested how subjects adapted to a change in the rewards associated with images. For  
32 this purpose, we instructed and trained subjects on new image-reward associations (see **State-**  
33 **space design**). Subjects then played two additional state-space blocks with these modified  
34 rewards.

35 Finally, we tested how subjects adapted to changes in the spatial structure of the state space.  
36 For this purpose, we told subjects that two pairs of images would switch locations, informing  
37 them precisely which images these were (see **State-space design**). Multiple-choice questions  
38 were used to ensure that subjects understood these instructions. Subjects then played a final  
39 state-space block with this modified spatial map.

40 At the end of the experiment, we also tested subjects’ explicit knowledge, asking them to  
41 sketch maps of the state spaces and indicate how many points each image was associated  
42 with before, and after, the reward contingency changed.

43 **Stimuli.** To ensure robust decoding from MEG, we used 8 images that differed in colour,  
44 shape, texture and semantic category<sup>43–45</sup>. These included: a frog, a face, a traffic sign, a  
45 tomato, a hand, a house, a pond, and a wrench.



1 **State-space task.** Subjects started each trial in a pseudorandom state, identified only by its  
2 associated image. Subjects then chose whether to move right, left, up, or down, and the  
3 chosen move was implemented on the screen, revealing the new state (i.e., as its associated  
4 image) to which the move led. In ‘one-move’ trials, this marked the end of the trial, and was  
5 followed by a short inter-trial interval. The next trial then started from another pseudorandom  
6 location. In ‘two-move’ trials, subjects made an additional move from the location where  
7 their first move had led. This second move disallowed backtracking the first move (e.g.,  
8 moving right and then left). Subjects were informed they would be awarded points associated  
9 with any image to which they move. Thus, subjects won points associated with a single  
10 image on one-move trials, and the combined value of the two images on two-move trials. The  
11 numbers of points awarded were never displayed during the main task. Every 6 trials, short  
12 text messages informed subjects what proportion of obtainable reward they had collected in  
13 the last 6 trials (message duration 2500 ms).

14 Each state-space block consisted 54 trials, 18 one-move and 36 two-move trials respectively.  
15 The first 6 trials were one-move, the next 12 were two-move trials, then the next 6 were again  
16 one-move trials, the next 12 two-move, and so on. Every 6 trials, short text messages  
17 informed subjects whether the next 6 trials were going to be one-move or two-move trials  
18 (message duration 2000 ms). Every six consecutive trials featured 6 different starting  
19 locations. The one exception to this were the first of the 24 two-move trials of the  
20 experiment, where in order to facilitate learning, each starting location repeated for two  
21 consecutive trials (a similar measure was also implemented for one-move trials during  
22 training; see **State-space training**). Subjects’ performance improved substantially in the  
23 second of such pairs of trials ( $\Delta$ proportion of optimal first choices = +0.15, 95% CI = +0.11  
24 to +0.18,  $p < 0.001$ , Bootstrap test).

25 At the beginning of every block (except the first one), we tested how well subjects could do  
26 the task without additional information, based solely on the identity of the starting locations.  
27 For this purpose, images to which subjects’ moves led were not shown for the first 12 trials.  
28 In two-move trials, this meant subjects implemented a second move from an unrevealed  
29 image (i.e., state).

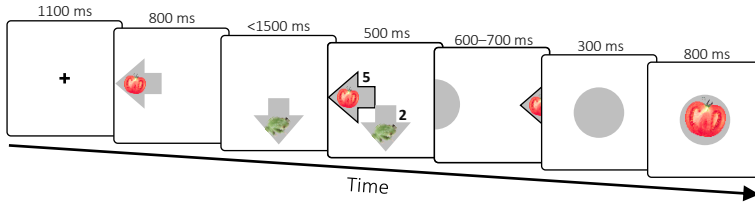
30 **State-space design.** The mapping of individual images to locations and rewards was  
31 randomly determined for each subject, but rewards were spatially organized in a similar  
32 manner for all subjects. To test whether subjects could flexibly adjust their choices, the state  
33 space was constructed such that there were five locations from which the optimal initial move  
34 was different depending on whether one or two moves were allowed. We tested subjects  
35 predominantly on these starting locations, using all five of them in every six consecutive  
36 trials. Following two blocks, the rewards associated with each image were changed, such that  
37 the optimal first moves in both 1-move and 2-move trials, given the new reward associations,  
38 were different from the optimal moves under the initial reward associations. The initial and  
39 modified reward associations were weakly anti-correlated across images ( $r = -0.37$ ).  
40 Finally, before the last block, we switched the locations of two pairs of images, such that the  
41 optimal first move changed for 15 out of 16 trial types (1- and 2-move trials x 8 starting  
42 locations).

43 **State-space training.** Subjects played six short training blocks, each block consisted 12 one-  
44 move trials starting in one of two possible locations. If a subject failed to collect 70% of the  
45 points available in one of these short blocks, the block was repeated. The majority of subjects  
46 (35 out of 40) had to repeat the first block, whereas only 12% of the remaining blocks were  
47 repeated (mean 0.6 blocks per subject, range 0 to 2). Very rarely, a block had to be repeated  
48 twice (a total of 5 out of 240 blocks for the whole group). Lastly, subjects played a final

1 training block consisting 48 one-move trials starting at any of the 8 possible locations. To  
2 facilitate learning, during the first half of the block, each starting location was repeated for  
3 two consecutive trials. In the second half of the block, starting locations were fully  
4 interleaved.

5 **State-space exercise.** Following the state-space training, which only included one-move  
6 trials, we ensured subjects understood how choices should differ in one- and two-move trials  
7 by asking them to choose the optimal moves in a series of random, fully visible state spaces.  
8 Subjects were given a bird's eye view of each state space, with each location showing the  
9 number of reward points with which it was associated. The starting location was indicated in  
10 addition to whether one, or two, moves were available from which to collect reward. In all  
11 exercises, the optimal initial move was different depending on whether one or two moves  
12 were allowed. Every 10 consecutive exercises consisted of 5 one-move trials and 5 two-move  
13 trials. To illustrate the continuity of the state space, the exercise included one-move and two-  
14 move trials, wherein the optimal move required the subject to move off the map and arrive at  
15 the other end (e.g., moving left from a leftmost location to arrive at the rightmost location). In  
16 another two-move trial, the optimal moves involved moving twice up or twice down, thereby  
17 returning to the starting location. Subjects continued to do the exercises until fulfilling a  
18 performance criterion of 9 correct answers in 10 consecutive exercises. This criterion was  
19 relaxed to 8 correct answers if at least 60 exercises had been completed. Only one subject  
20 required 60 exercises to reach criterion (mean required exercises = 24.5 exercises, SD 9.3).

21 **Image-reward training.** To ensure subjects remembered how many points each image  
22 awarded, we required subjects to select the more rewarding image out of any pair of  
23 presented images. First, subjects were asked to memorize the number of points each image  
24 would award. Then, each round of training consisted of 28 trials, testing subjects on all 28  
25 possible pairs of images (**Supplementary Fig. 5**). Each trial started with the presentation of  
26 one image, depicted on an arrow pointing either right, left, up or down. 800 ms later, another  
27 image appeared on an arrow pointing in a different direction. Subjects had then to press the  
28 button corresponding to the direction of the more rewarding image. Here, as throughout the  
29 experiment, subjects were instructed to press the 'left' and 'up' buttons with their left hand,  
30 and the 'right' and 'down' buttons with their right hand. During training, images were  
31 mapped to directions such that each of the four directions was equally associated with low-  
32 and high-reward images. Once subjects made their choice, the number of points associated  
33 with each of the two images appeared on the screen, and if the choice was correct the chosen  
34 move was implemented on the screen. Subjects repeated this training until they satisfied a  
35 performance criterion, based on how many points they missed consequent upon choosing less  
36 rewarding images. The initial performance criterion allowed 4 missed points, or less, in a  
37 whole training round (out of a maximum of 130 points). This criterion was gradually relaxed,  
38 to 8 missed points in the second training round, to 12 missed points in the third training  
39 round, and to 16 missed points thereafter. Once subjects satisfied the performance criterion  
40 without time limit, they repeated the training with only 1500 ms allowed to make each  
41 choice, until satisfying the same re-set gradually relaxing criterion. Overall, subjects required  
42 an average of 3.4 training rounds (SD 1.0) to learn the initial image-reward associations (1.3  
43 rounds without, and then 2.1 rounds with, a 1500 ms time limit), and 4.3 rounds (SD 1.3) to  
44 learn the second set image-reward associations (2.0 rounds without, and 2.3 rounds with, a  
45 time limit). Questioning at the end of the experiment validated that subjects had explicit  
46 recall for both sets of image-reward associations (mean error 0.36 pts, SEM = 0.07 pts;  
47 chance = 4.05 pts).



**Supplementary Fig. 5.** Image-reward training. Timeline of a trial.

1 **Modelling.** To test what decision algorithm subjects employed, and in particular, whether  
2 they chose moves that had previously been most rewarding from the same starting location  
3 (model-free planning), or whether they learned how the state space is structured and used this  
4 information to plan ahead (model-based planning), we compared between model-free and  
5 model-based algorithms in terms of how well they fitted subjects' actual choices. These  
6 models were informed by previous work<sup>46,47</sup>, adjusted to the present task, and validated using  
7 model and parameter recovery tests on simulated data.

8 **Model-free learning algorithm** (free parameters:  $\eta^{\text{MF1}}$ ,  $\eta^{\text{MF2}}$ ,  $\tau^{\text{MF}}$ ,  $\tau'^{\text{MF}}$ ,  $\theta$ ,  $\beta_{1,2}^{\text{MF1}}$ ,  $\beta_2^{\text{MF2}}$ ,  
9  $\gamma_{\text{up,down,left,right}}$ ). This algorithm learns the expected value of performing a given move upon  
10 encountering a given image. To do this, the algorithm updates its expectation  $Q^{\text{MF}}$  from move  
11  $m$  given image  $s$  whenever this move is taken and its outcome is observed:

$$12 \quad Q_{t+1}^{\text{MF1}}(s_{t,1}, m_t) = Q_t^{\text{MF1}}(s_{t,1}, m_t) + \eta^{\text{MF1}} \delta_t^{\text{MF1}}, \quad (1)$$

13 where  $s_{t,1}$  is trial  $t$ 's starting image,  $\delta_t^{\text{MF}}$  is the reward prediction error, and  $\eta^{\text{MF1}}$  is a fixed  
14 learning rate between 0 and 1. Reward prediction errors are computed as the difference  
15 between actual and expected outcomes:

$$16 \quad \delta_t^{\text{MF1}} = R_g(s_{t,2}) - Q_t^{\text{MF1}}(s_{t,1}, m_t), \quad (2)$$

17 where the actual outcome consists of the points associated with the new image to which the  
18 move led,  $R_g(s_{t,2})$ .  $g = 1$  refers to the initial image-rewards associations, and  $g = 2$  refers  
19 to the second set of image-rewards associations about which subjects were instructed in the  
20 middle of the experiment.

21 On 2-move trials, the algorithm also learns the expected reward for each pair of moves given  
22 each starting image. Thus, another set of Q values is maintained ( $Q^{\text{MF2}}$ ), one for each  
23 possible pair of moves for each starting image, and these are updated every time a pair of  
24 moves is completed based on the total reward obtained by the two moves. This learning  
25 proceeds as described by Eqs. 1 and 2, but with a different learning rate ( $\eta^{\text{MF2}}$ ).

26 All expected values are initialized to  $\theta$ , and decay back to this initial value before every  
27 update:

$$28 \quad Q^{\text{MF}} \leftarrow \tau^{\text{MF}} Q^{\text{MF}} + (1 - \tau^{\text{MF}}) \theta, \quad (3)$$

29 where  $\tau^{\text{MF}}$  value retention. This allows learned expectations to be gradually forgotten.

30 Following instructed changes to the number of points associated with each image, or to the  
31 spatial arrangement of the images, previously learned Q values are of little use. Thus, we  
32 allow the Q values to then return back to  $\theta$ , as in Eq. 3, but only for a single timestep and  
33 with a different, potentially lower, memory parameter  $\tau'^{\text{MF}}$ .

- 1 Finally, the algorithm chooses moves based on a combination of its learned expected values.  
 2 On 1-move trials, only single-move Q values are considered:

$$3 \quad p(m_t = a|s_t) \propto e^{\gamma_m + \beta_1^{\text{MF1}} Q_t^{\text{MF1}}(s_{t,1}, m)}, \quad (4)$$

4 where  $\gamma_m$  is a fixed bias in favor of move  $m$  ( $\sum_m \gamma_m = 0$ ), and  $\beta_1^{\text{MF1}}$  is an inverse  
 5 temperature parameter that weighs the impact of expected values on choice. On 2-move  
 6 trials, both types of Q values are considered. Thus, the first move is chosen based on a  
 7 weighted sum of the single-move Q values and the move-pair Q values:

$$8 \quad p(m_{t,1} = m|s_{t,1}) \propto e^{\gamma_m + \beta_2^{\text{MF1}} Q_t^{\text{MF1}}(s_{t,1}, m) + \beta_2^{\text{MF2}} Q_t^{\text{MF2}}(s_{t,1}, m)}, \quad (5)$$

9 wherein the latter are integrated over possible second moves each weighted by its probability:

$$10 \quad Q_t^{\text{MF2}}(s_{t,1}, m) = \sum_{m^*} p(m_{t,2} = m^*|s_{t,1}, m_{t,1}) Q_t^{\text{MF2}}(s_{t,1}, m, m^*) \quad (6)$$

11 Then, in choosing the second move the algorithm takes into account the state to which the  
 12 first move led:

$$13 \quad p(m_{t,2} = m|s_{t,1}, m_{t,1}, s_{t,2}) \propto e^{\gamma_m + \beta_2^{\text{MF1}} Q_t^{\text{MF1}}(s_{t,2}, m) + \beta_2^{\text{MF2}} Q_t^{\text{MF2}}(s_{t,1}, m_{t,1}, m)}. \quad (7)$$

14 However, when the newly reached image  $s_{t,2}$  is not known (i.e., in trials without feedback, or  
 15 when estimating  $p(m_{t,2} = m^*|s_{t,1}, m_{t,1})$  in Eq. 6 before  $s_{t,2}$  is reached),  $Q^{\text{MF1}}$  values are  
 16 averaged over all settings of  $s_{t,2}$ .

17 **Model-based learning algorithm** (free parameters:  $\eta^{\text{MB}}, \tau^{\text{MB}}, \tau'^{\text{MB}}, \rho, \omega, \beta^{\text{MB}}, \kappa,$   
 18  $\gamma_{\text{up,down,left,right}}$ ). This algorithm learns the probability of transitioning from one image to  
 19 another following each move. To do this, the algorithm updates its probability estimates,  $T$ ,  
 20 whenever a move is made and a transition is observed:

$$21 \quad T_{t+1}(s_{t,1}, m_t, s_{t,2}) = T_t(s_{t,1}, m_t, s_{t,2}) + \eta^{\text{MB}} \delta_t^{\text{MB}}, \quad (8)$$

22 where  $\delta_t^{\text{MB}}$  is the image-transition prediction error, and  $\eta^{\text{MF}}$  is a fixed learning rate between 0  
 23 and 1. Image-transition prediction errors reflect the difference between actual and expected  
 24 transitions:

$$25 \quad \delta_t^{\text{MB}} = 1 - T_t(s_{t,1}, m_t, s_{t,2}). \quad (9)$$

26 To ensure that transition probabilities sum to 1, the transition matrix is renormalized  
 27 following every update:

$$28 \quad \forall s \quad T_{t+1}(s_{t,1}, m_t, s) \leftarrow \frac{T_{t+1}(s_{t,1}, m_t, s)}{\sum_{s'} T_{t+1}(s_{t,1}, m_t, s')}. \quad (10)$$

29 Learning may also take place with respect to the opposite transition. For instance, if moving  
 30 right from image  $s_{t,1}$  leads to image  $s_{t,2}$ , the agent can infer that moving left from image  $s_{t,2}$   
 31 would lead to image  $s_{t,1}$ . Such inference is modulated in the algorithm by free parameter  $\rho$ :

$$32 \quad T_{t+1}(s_{t,1}, \tilde{m}_t, s_{t,2}) = T_t(s_{t,1}, \tilde{m}_t, s_{t,2}) + \rho \eta^{\text{MB}} \delta'_t{}^{\text{MB}}, \quad (11)$$

33 where  $\tilde{m}_t$  is the opposite of  $m_t$ , and  $\delta'$  is the opposite transition prediction error:

$$34 \quad \delta'_t{}^{\text{MB}} = 1 - T_t(s_{t,2}, \tilde{m}_t, s_{t,1}). \quad (12)$$

1 Self-transitions are impossible and thus their probability is initialized to 0. All other  
 2 transitions are initialized with uniform probabilities, and these probabilities decay back to  
 3 their initial values before every update:

$$4 \quad T \leftarrow \tau^{\text{MB}}T + (1 - \tau^{\text{MB}})\frac{1}{7}, \quad (13)$$

5 where  $\tau^{\text{MB}}$  is the model-based memory parameter. A low  $\tau^{\text{MB}}$  results in faster decay of  
 6 expected transition probabilities towards uniform distributions, decreasing the impact of MB  
 7 knowledge on choice.

8 When instructed about changes to the image locations, the agent rearranges its transition  
 9 probabilities based on the instructed changes with limited success, as indexed by free  
 10 parameter  $\omega$ :

$$11 \quad T \leftarrow (1 - \omega)T + \omega T^{\text{rearranged}}. \quad (14)$$

12 Since some subjects may simply reset their transition matrix following instructed changes,  
 13 the algorithm also ‘forgets’ after such instruction, as in Eq. 13, but only for a single time  
 14 point and with a different memory parameter,  $\tau'^{\text{MB}}$ .

15 Finally, the probability the algorithm will choose a given move when encountering a given  
 16 image depends on its model-based estimate of the move’s expected outcome:

$$17 \quad p(m_t = m | s_{t,1}) \propto e^{\gamma m + \beta^{\text{MB}} Q_t^{\text{MB}}(s_{t,1}, m)}. \quad (15)$$

18 The algorithm estimates expected outcomes by multiplying the number of points associated  
 19 with an image with the probability of transitioning to that image, integrating over all potential  
 20 future images:

$$21 \quad Q_t^{\text{MB}}(s_{t,1}, m) = \sum_s T_t(s_{t,1}, m, s) R_g(s). \quad (16)$$

22 When two moves are allowed, the calculation also accounts for the number of points  
 23 obtainable with the second move,  $m_{t,2}$ :

$$24 \quad Q_t^{\text{MB}}(s_{t,1}, m) = \sum_s T_t(s_{t,1}, m, s) \left( R_g(s) + \kappa \max_{m'} \sum_{s'} T_t(s, m', s') R_g(s') \right), \quad (17)$$

25 where  $\kappa$  is a fractional parameter that determines the degree to which reward obtained by the  
 26 second move is taken into account.

27 Following the first move, Eq. 15 is used to choose a second move based on the observed new  
 28 location ( $s_{t,2}$ ). However, if the next location is not shown (i.e., in trials without feedback), the  
 29 agent chooses its second move by integrating Eq. 15 over the expected  $s_{t,2}$ , as determined by  
 30  $T_t(s_{t,1}, m_{t,1}, s_{t,2})$ .

31 **MF-MB hybrid algorithm.** This algorithm employs both model-free (MF) and model-based  
 32 (MB) planning, choosing moves based on a combination of the expected values estimated by  
 33 the two learning processes:

$$34 \quad p(m_t = m | s_t) \propto e^{\gamma m + \beta_1^{\text{MF}} Q_t^{\text{MF}}(s_{t,1}, m) + \beta^{\text{MB}} Q_t^{\text{MB}}(s_{t,1}, m)}, \quad (18)$$

1 In 2-move trials, the algorithm makes a choice based on a combination of the model-based Q  
2 values and both the single-move and two-move model-free Q values. For the first move, the  
3 combination is:

$$4 \quad p(m_{t,1} = m | s_{t,1}) \propto e^{\gamma_m + \beta_2^{\text{MF1}} Q_t^{\text{MF1}}(s_{t,1}, m) + \beta_2^{\text{MF2}} Q_t^{\text{MF2}}(s_{t,1}, m) + \beta^{\text{MB}} Q^{\text{MB}}(s_{t,1}, m)}, \quad (19)$$

5 with  $Q^{\text{MB}}(s_{t,1}, m)$  computed according to Eq. 17. For the second move, the choice is made  
6 according to:

$$7 \quad p(m_{t,1} = m | s_{t,2}) \propto e^{\gamma_m + \beta_2^{\text{MF1}} Q_t^{\text{MF1}}(s_{t,2}, m) + \beta_2^{\text{MF2}} Q_t^{\text{MF2}}(s_{t,1}, m_{t,1}, m) + \beta^{\text{MB}} Q^{\text{MB}}(s_{t,2}, m)}. \quad (20)$$

8 When the image is not shown following the first move (i.e., in a no-feedback trial), the agent  
9 averages the model-free values over all images.

10 **Parameter fitting.** To fit the free parameters of the different algorithms to subjects' choices,  
11 we used an iterative hierarchical expectation-maximization procedure<sup>26</sup>. We first sampled  
12 10000 random settings of the parameters from predefined group-level prior distributions.  
13 Then, we computed the likelihood of observing subjects' choices given each setting, and used  
14 the computed likelihoods as importance weights to re-fit the parameters of the group-level  
15 prior distributions. These steps were repeated iteratively until model evidence ceased to  
16 increase (see **Model Comparison** below for how model evidence was estimated). This  
17 procedure was then repeated with 31623 samples per iteration, and finally with 100000  
18 samples per iteration. To derive the best-fitting parameters for each individual subject, we  
19 computed a weighted mean of the final batch of parameter settings, in which each setting was  
20 weighted by the likelihood it assigned to the subject's choices. Fractional parameters ( $\eta^{\text{MF}}$ ,  
21  $\tau^{\text{MF}}$ ,  $\tau'^{\text{MF}}$ ,  $\eta^{\text{MB}}$ ,  $\tau^{\text{MB}}$ ,  $\tau'^{\text{MB}}$ ,  $\rho$ ,  $\omega$ ,  $\alpha$ ) were modelled with Beta distributions (initialized with  
22 shape parameters  $a = 1$  and  $b = 1$ ) and their values were log-transformed for the purpose of  
23 subsequent analysis. Initial Q values ( $\theta$ ) and bias parameters ( $\gamma_{\text{up}}$ ,  $\gamma_{\text{down}}$ ,  $\gamma_{\text{left}}$ ,  $\gamma_{\text{right}}$ ) were  
24 modelled with normal distributions (initialized with  $\mu = 0$  and  $\sigma = 1$ ) to allow for both  
25 positive and negative effects, and all other parameters were modeled with Gamma  
26 distributions (initialized with shape = 1, scale = 1).

27 **Algorithm comparison.** We compared between pairs of algorithms, in terms of how well  
28 each accounted for subjects' choices, by means of the integrated Bayesian Information  
29 Criterion (iBIC)<sup>48,49</sup>. To do this, we estimated the evidence in favour of each model ( $\mathcal{L}$ ) as the  
30 mean likelihood of the model given 100000 random parameter settings drawn from the fitted  
31 group-level priors. We then computed the iBIC by penalizing the model evidence to account  
32 for algorithm complexity as follows:  $\text{iBIC} = -2 \ln \mathcal{L} + k \ln n$ , where  $k$  is the number of fitted  
33 parameters and  $n$  is the number of subject choices used to compute the likelihood. Lower  
34 iBIC values indicate a more parsimonious fit.

35 **Algorithm and parameter recovery tests.** We tested whether our dataset was sufficiently  
36 informative to distinguish between the MF, MB and hybrid algorithms and recover the  
37 correct parameter values. For this purpose, we generated 10 simulated datasets using each  
38 algorithm and applied our fitting and comparison procedures to each dataset. To reduce  
39 processing time, only 10000 parameter settings were sampled. To maximize the chances of  
40 confusion between algorithms, we implemented all algorithms with the parameter values that  
41 best fitted subjects' choices. Algorithm comparison implicated the correct algorithm in each  
42 of the 30 simulated datasets, and the parameters values that best fitted the simulated data  
43 consistently correlated with the actual parameter values used to generate these data  
44 (Pearson's  $r$ :  $M = 0.57$ ,  $SEM = 0.05$ ), and this correlation was stronger for parameters

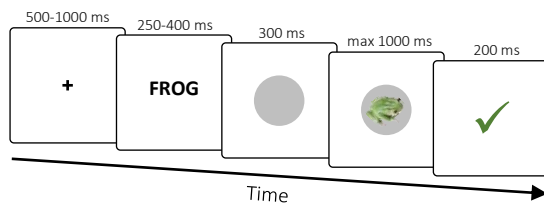
1 whose values were used for multiple trials when computing the fit to data (e.g., learning rates  
2 and inverse temperature parameters;  $M = 0.67$ ,  $SEM = 0.04$ ).

3 **Additional algorithms.** To test whether the algorithms described above were most suitable  
4 for describing subjects' behaviour, we compared them to several additional algorithms, all of  
5 which failed to fit subjects' choices as well as the above counterparts, and so we do not  
6 describe them in detail. These alternative algorithms included a MF algorithm that only learns  
7 single-move Q values, but employs temporal difference learning<sup>50</sup> to backpropagate second  
8 outcomes in 2-move trials back to the Q values of the starting location (BIC = 41559); a MB  
9 algorithm that employs Bayesian inference with a uniform Dirichlet prior<sup>26</sup> to learn the  
10 multinomial distributions that compose the state transition matrix (BIC = 43301); a MF-MB  
11 hybrid algorithm where state-transition expectations are only used to account for prospective  
12 second-move Q values when choosing the first move in 2-move trials (BIC = 40920); and an  
13 algorithm that combines two MF algorithms with different parameters (BIC = 40715).

14 **MEG acquisition.** MEG was recorded continuously at 600 samples/second using a whole-  
15 head 275-channel axial gradiometer system (CTF Omega, VSM MedTech, Canada), while  
16 subjects sat upright inside the scanner. A projector displayed the task on a screen ~80 cm in  
17 front of the subject. Subjects made responses by pressing a button box, using their left hand  
18 for 'left' and 'up' choices and their right hand for 'right' and 'down' choices. Pupil size and  
19 eye gaze were recorded at 250 Hz using a desktop-mounted EyeLink II eyetracker (SR  
20 Research).

21 **MEG preprocessing.** Preprocessing was performed using the Fieldtrip toolbox<sup>51</sup> in  
22 MATLAB (MathWorks). Data from two sensors were not recorded due to a high level of  
23 noise detected in routine testing. Data were first manually inspected for jump artefacts. Then,  
24 independent component analysis was used to remove components that corresponded to eye  
25 blinks, eye movement and heart beats. Based on previous experience<sup>25</sup>, we expected stimuli  
26 to be represented in low frequency fluctuations of the MEG signal. Therefore, to remove fast  
27 muscle artefacts and slow movement artefacts, we low-pass filtered the data with a 20 Hz  
28 cutoff frequency using a sixth-order Butterworth IIR filter, and we baseline-corrected each  
29 trial's data by subtracting the mean signal recorded during the 400 ms preceding trial onset.  
30 Trials in which the average standard deviation of the signal across channels was at least 3  
31 times greater than median were excluded from analysis (0.4% of trials, SEM 0.2%). Finally,  
32 the data were resampled from 600 Hz to 100 Hz to conserve processing time and improve  
33 signal to noise ratio. Therefore, data samples used for analysis were length 273 vectors  
34 spaced every 10 ms.

35 **Pre-task stimulus exposure.** To allow decoding of images from MEG we instructed subjects  
36 to identify each of the images in turn (**Supplementary Fig. 6**). On each trial, the target image  
37 was indicated textually (e.g., 'FACE') and then an image appeared on the screen. Subjects'  
38 task was to report whether the image matched (LEFT button) or did not match (RIGHT  
39 button) the preceding text. 20% of presented images did not match the text. The task  
40 continued until subjects correctly identified each of the images at least 25 times. Subjects  
41 were highly accurate on both match ( $M = 97.2\%$ ,  $SEM = 0.4\%$ ) and no-match ( $M = 90.2\%$ ,  
42  $SEM = 0.6\%$ ) trials. To ensure robust decoding from MEG, we chose eight images that  
43 differed in colour, shape, texture and semantic category<sup>43,44</sup> (**Fig. 1a**). Importantly, at this  
44 point subjects had no knowledge as to what the main task would involve, nor that the images  
45 would be associated with state-space locations and rewards. This ensured that no task  
46 information could be represented in the MEG data at this stage.



**Supplementary Fig. 6.** Pre-task stimulus exposure. Timeline of a trial.

1 **MEG decoding.** We used support vector machines (SVMs) to decode images and moves  
2 from MEG. All decoders were trained on MEG data recorded outside of the main state-space  
3 task and validated within the task. As in previous work<sup>25</sup>, we trained a separate decoder for  
4 each time bin between 150 and 600 ms following the relevant event, either image onset or  
5 move choice, resulting in 46 decoders whose output was averaged. Averaging over decoders  
6 trained at different time points reduces peak decodability following stimulus onset, but can  
7 increase decodability of stimuli that are being processed when not on the screen<sup>25</sup>. To avoid  
8 over-fitting, training and testing were performed on separate sets of trials following a 5-fold  
9 cross validation scheme. These analyses were performed using LIBSVM's implementation of  
10 the C-SVC algorithm with radial basis functions<sup>52</sup>. Decoder training and testing were  
11 performed with each of 16 combinations of the algorithms' cost parameter ( $10^{-1}$ ,  $10^0$ ,  $10^1$ ,  
12  $10^2$ ) and basis-function concentration parameter ( $10^{-2}/n$ ,  $10^{-1}/n$ ,  $10^0/n$ ,  $10^1/n$ ), where  $n$   
13 is the number of MEG features (273 channels). Where classes differed in number of  
14 instances, weighting was used to ensure classes were equally weighted.

15 To decode the probability of each of eight possible images being presented (8-way  
16 classification), we used MEG data recorded during pre-task stimulus exposure. Decoding was  
17 evaluated based on the mean probability the decoders assigned to the presented image. To  
18 decode the probability of each of the four possible moves (LEFT, RIGHT, UP, DOWN)  
19 being chosen (4-way classification), we used MEG data recorded during the image-reward  
20 training. For both types of decoder, the parameter combination of cost =  $10^2$  and  
21 concentration =  $10^{-2}/n$  yielded the best cross-validated decoding performance and was thus  
22 used for all ensuing analyses.

23 **Sequenceness measure.** To investigate how representations of different images related to  
24 one another in time, we used a measure recently developed for detecting sequences of  
25 representations in MEG<sup>10</sup>. 'Sequenceness' is computed as the difference between the cross-  
26 correlation of two images' decodability time-series with positive and negative time lags. By  
27 relying on asymmetries in the cross-correlation function, this measure detects sequential  
28 relationships even between closely correlated (or anti-correlated) time series, as we have  
29 previously demonstrated on simulated time series<sup>25</sup>. Positive values indicate that changes in  
30 the first time series are followed by similar changes in the second time series ('forward  
31 sequenceness'), negative values indicate the reverse sequence ('backward sequenceness'),  
32 and zero indicates no sequential relationship. As in previous work, cross correlations were  
33 computed between the z-scored time series over 400 ms sliding windows with time lags of up  
34 to 200 ms. This timescale is sufficient for capturing the relationship between successive alpha  
35 cycles, which is important given the possibility that such oscillations may reflect temporal  
36 quanta of information processing<sup>53</sup>.

37 **Bayesian hierarchical Gaussian Process time series analysis.** To determine whether  
38 sequenceness time-series recorded following outcomes provided robust evidence of replay  
39 that correlated with individual index of flexibility, we modelled each mean sequenceness  
40 time-series as a summation of two zero-mean Gaussian Processes with squared exponential  
41 kernels: a group-level process and an individual-level process. The group-level process



1 identifies the timepoints in which sequenceness systematically deviates from zeros, and the  
2 individual-level processes (one for each time series) account for deviations of individual time  
3 series from the group-level process.

4 To enable completion of the MCMC sampling within a reasonable timeframe, we reduced the  
5 trial-to-trial sequenceness data to four mean time series per subject: sequenceness encoding  
6 the last or penultimate transition following highly or weakly surprising outcomes. High and  
7 low surprise were determined based on the state prediction error generated by the hybrid  
8 algorithm, whose parameters were fitted to the individual subject's choices (i.e., high –  
9 above-mean prediction error, low – below-mean prediction error). Since we assumed last and  
10 penultimate transitions could be replayed in different timepoints, these two types of time  
11 series each had their own group-level Gaussian Process. To account for the factors of IF and  
12 surprise, the group-level process was multiplied for each time series by a weighted linear  
13 combination of the two factors, their interaction, and an intercept (thus involving four  
14 parameters:  $\beta$ ,  $\beta^{\text{subject}}$ ,  $\beta^{\text{surprise}}$ ,  $\beta^{\text{interaction}}$ ). The two types of Gaussian Process were  
15 parameterized by different length-scales ( $\rho^{\text{group}}$ ,  $\rho^{\text{individual}}$ ) and marginal standard  
16 deviations ( $\alpha^{\text{group}}$ ,  $\alpha^{\text{individual}}$ ), and an a standard deviation parameter ( $\sigma$ ) accounted for  
17 additional normally distributed noise across all observations.

18 Bayesian estimation was performed in R<sup>54</sup> using the STAN55 package for Markov Chain  
19 Monte Carlo (MCMC) sampling. Prior distributions were set so as to be weakly informative  
20 and have broad range on the scale of the variables<sup>29</sup>. Thus,  $\beta$  coefficients were drawn from  
21 normal distributions with a mean of zero and a standard deviation of 10. All predictor  
22 variables were standardized. Standard deviations parameters ( $\alpha$ ,  $\sigma$ ) were drawn from a  
23 truncated normal distribution limited to positive values, with a mean of zero and a standard  
24 deviation that matches the standard deviation of the predicted variable. Length-scales ( $\rho$ )  
25 were drawn from log-normal distributions whose mean is the geometric mean of two  
26 extremes: the distance in time between two successive timepoints, and the distance in time  
27 between the first and last timepoints. Half of the difference between these two values was  
28 used as the standard deviation of the priors.  $\beta^{\text{interaction}}$  was limited to positive values for the  
29 sake of identifiability, since the group-level Gaussian Processes were multiplied by the  $\beta$   
30 coefficients.

31 We ran six MCMC chains each for 1400 iterations, with the initial 400 samples used for  
32 warmup. STAN's default settings were used for all other settings. Examining the results  
33 showed there were no divergent transitions, and all parameters were estimated with effective  
34 sample sizes larger than 1000 and shrink factors smaller than 1.1. Posterior predictive checks  
35 showed good correspondence between the real and generated data (**Fig. 2b,c**).

36 **Decodability time series analyses.** Decodability was tested for difference from zero and  
37 covariance with individual flexibility using the Bayesian Gaussian Process approach outlined  
38 above with the exclusion of the surprise predictor, which is inapplicable to timepoints that  
39 precede outcome onset.

40 **Other statistical Methods.** Significance tests were conducted using nonparameteric methods  
41 that do not assume specific distributions. Differences from zero were tested using 10000  
42 samples of bias-corrected and accelerated Bootstrap with default MATLAB settings.  
43 Correlations and differences between groups were tested by comparison to null distributions  
44 generated by 10000 permutations of the pairing between the two variables of interest. All  
45 tests are two-tailed.

1 **Data and Code availability**

2 The data and custom code used in this study have been deposited on the Open Science  
3 Framework under DOI 10.17605/OSF.IO/GUHJE.

4 **Acknowledgements**

5 We thank Zeb Kurth-Nelson and Yunzhe Liu for helpful comments on a previous version of  
6 the manuscript. E.E. holds an Alon Fellowship from the Israeli Council for Higher Education.  
7 P.D. is funded by the the Max Planck Society. R.J.D. holds a Wellcome Trust Investigator  
8 award (098362/Z/12/Z). The Max Planck UCL Centre for Computational Psychiatry and  
9 Ageing Research is a joint initiative supported by the Max Planck Society and University  
10 College London. The Wellcome Centre for Human Neuroimaging is supported by core  
11 funding from the Wellcome Trust (091593/Z/10/Z).

12 **Competing interests**

13 The authors declare that they have no conflict of interest.

## References

1. Olafsdottir, H. F., Carpenter, F., & Barry, C. (2017). Task demands predict a dynamic switch in the content of awake hippocampal replay. *Neuron*, 96(4), 925-935.
2. Pezzulo, G., van der Meer, M. A., Lansink, C. S., & Pennartz, C. M. (2014). Internally generated sequences in learning and executing goal-directed behavior. *Trends in cognitive sciences*, 18(12), 647-657.
3. Diba, K., & Buzsáki, G. (2007). Forward and reverse hippocampal place-cell sequences during ripples. *Nature neuroscience*, 10(10), 1241.
4. Foster, D. J., & Wilson, M. A. (2006). Reverse replay of behavioural sequences in hippocampal place cells during the awake state. *Nature*, 440(7084), 680.
5. Louie, K., & Wilson, M. A. (2001). Temporally structured replay of awake hippocampal ensemble activity during rapid eye movement sleep. *Neuron*, 29(1), 145-156.
6. Skaggs, W. E., & McNaughton, B. L. (1996). Replay of neuronal firing sequences in rat hippocampus during sleep following spatial experience. *Science*, 271(5257), 1870-1873.
7. Ji, D., & Wilson, M. A. (2007). Coordinated memory replay in the visual cortex and hippocampus during sleep. *Nature neuroscience*, 10(1), 100.
8. Gupta, A. S., van der Meer, M. A., Touretzky, D. S., & Redish, A. D. (2010). Hippocampal replay is not a simple function of experience. *Neuron*, 65(5), 695-705.
9. Ólafsdóttir, H. F., Barry, C., Saleem, A. B., Hassabis, D., & Spiers, H. J. (2015). Hippocampal place cells construct reward related sequences through unexplored space. *Elife*, 4, e06063.
10. Kurth-Nelson, Z., Economides, M., Dolan, R. J., & Dayan, P. (2016). Fast sequences of non-spatial state representations in humans. *Neuron*, 91(1), 194-204.
11. Schuck, N. W., & Niv, Y. (2019). Sequential replay of nonspatial task states in the human hippocampus. *Science*, 364(6447), eaaw5181.
12. Liu, Y., Dolan, R. J., Kurth-Nelson, Z., & Behrens, T. E. (2019). Human Replay Spontaneously Reorganizes Experience. *Cell*.
13. Foster, D. J. (2017). Replay comes of age. *Annual review of neuroscience*, 40, 581-602.
14. Behrens, T. E., Muller, T. H., Whittington, J. C., Mark, S., Baram, A. B., Stachenfeld, K. L., & Kurth-Nelson, Z. (2018). What is a cognitive map? Organizing knowledge for flexible behavior. *Neuron*, 100(2), 490-509.
15. Stachenfeld, K. L., Botvinick, M. M., & Gershman, S. J. (2017). The hippocampus as a predictive map. *Nature neuroscience*, 20(11), 1643.
16. Momennejad, I., Otto, A. R., Daw, N. D., & Norman, K. A. (2018). Offline replay supports planning in human reinforcement learning. *Elife*, 7, e32548.
17. Mattar, M. G., & Daw, N. D. (2018). Prioritized memory access explains planning and hippocampal replay. *Nature Neuroscience*, 21(11), 1609.
18. Sutton, R. S. (1991). Dyna, an integrated architecture for learning, planning, and reacting. *ACM Sigart Bulletin*, 2(4), 160-163.
19. Gershman, S. J., Markman, A. B., & Otto, A. R. (2014). Retrospective revaluation in sequential decision making: A tale of two systems. *Journal of Experimental Psychology: General*, 143(1), 182.

20. Kurdi, B., Gershman, S. J., & Banaji, M. R. (2019). Model-free and model-based learning processes in the updating of explicit and implicit evaluations. *Proceedings of the National Academy of Sciences*, *116*(13), 6035-6044.
21. Crockett, M. J. (2013). Models of morality. *Trends in cognitive sciences*, *17*(8), 363-366.
22. Everitt, B. J., & Robbins, T. W. (2005). Neural systems of reinforcement for drug addiction: from actions to habits to compulsion. *Nature neuroscience*, *8*(11), 1481.
23. Gillan, C. M., Fineberg, N. A., & Robbins, T. W. (2017). A trans-diagnostic perspective on obsessive-compulsive disorder. *Psychological medicine*, *47*(9), 1528-1548.
24. Daw, N. D., Gershman, S. J., Seymour, B., Dayan, P., & Dolan, R. J. Model-based influences on humans' choices and striatal prediction errors. *Neuron* *69*, 1204-1215 (2011).
25. Eldar, E., Bae, G. J., Kurth-Nelson, Z., Dayan, P., & Dolan, R. J. (2018). Magnetoencephalography decoding reveals structural differences within integrative decision processes. *Nature Human Behaviour*, *2*(9), 670.
26. Bishop, C.M. (2006) *Pattern Recognition and Machine Learning* (Springer).
27. Pfeiffer, B. E., & Foster, D. J. (2013). Hippocampal place-cell sequences depict future paths to remembered goals. *Nature*, *497*(7447), 74.
28. Kurth-Nelson, Z., Barnes, G., Sejdinovic, D., Dolan, R. & Dayan, P. Temporal structure in associative retrieval. *Elife* *4*, e04919 (2015).
29. Kruschke, J. (2014). *Doing Bayesian data analysis: A tutorial with R, JAGS, and Stan*. Academic Press.
30. Moore, A. W., Atkeson, C. G. (1993) Prioritized sweeping: Reinforcement learning with less data and less time. *Machine Learning*, *13*, 103–130.
31. Peng, J., & Williams, R. J. (1993) Efficient learning and planning within the Dyna framework. *IEEE International Conference on Neural Networks* 168–174. DOI: <https://doi.org/10.1109/ICNN.1993.298551>.
32. Kahneman, D. (2011). *Thinking, fast and slow*. Macmillan.
33. Stanovich, K. E., & West, R. F. (2000). Individual differences in reasoning: Implications for the rationality debate? *Behavioral and brain sciences* *23*, 645-665.
34. Gläscher, J., Daw, N., Dayan, P., & O'Doherty, J. P. (2010). States versus rewards: dissociable neural prediction error signals underlying model-based and model-free reinforcement learning. *Neuron* *66*, 585-595.
35. Decker, J. H., Otto, A. R., Daw, N. D., & Hartley, C. A. (2016). From creatures of habit to goal-directed learners: Tracking the developmental emergence of model-based reinforcement learning. *Psychological science* *27*, 848-858.
36. Gillan, C. M., Otto, A. R., Phelps, E. A., & Daw, N. D. (2015). Model-based learning protects against forming habits. *Cognitive, Affective, & Behavioral Neuroscience* *15*, 523-536.
37. da Silva, C. F., & Hare, T. (2019). Model-free or muddled models in the two-stage task? *bioRxiv* 682922.
38. Kool, W., Cushman, F. A., & Gershman, S. J. (2016). When does model-based control pay off? *PLoS computational biology* *12*, e1005090.

39. Akam, T., Rodrigues-Vaz, I., Zhang, X., Pereira, M., Oliveira, R., Dayan, P., & Costa, R. M. (2017). Single-Trial Inhibition of Anterior Cingulate Disrupts Model-based Reinforcement Learning in a Two-step Decision Task. *bioRxiv* 126292.
40. Russek, E. M., Momennejad, I., Botvinick, M. M., Gershman, S. J., & Daw, N. D. (2017). Predictive representations can link model-based reinforcement learning to model-free mechanisms. *PLoS computational biology*, *13*(9), e1005768.
41. Carey, A. A., Tanaka, Y., & Van Der Meer, M. (2019). Reward revaluation biases hippocampal replay content away from the preferred outcome. *Nature Neuroscience* *22*, 1450-1459 (2019).
42. Hunt, L. T. et al. Mechanisms underlying cortical activity during value-guided choice. *Nat. Neurosci.* *15*, 470–476 (2012)
43. Carlson, T., Tovar, D. A., Alink, A. & Kriegeskorte, N. Representational dynamics of object vision: the first 1000 ms. *J. vis.* *13*, 1 (2013)
44. Isik, L., Meyers, E. M., Leibo, J. Z. & Poggio, T. The dynamics of invariant object recognition in the human visual system. *J. Neurophysiol.* *111*, 91–102 (2014)
45. Cichy, R. M., Pantazis, D. & Oliva, A. Resolving human object recognition in space and time. *Nat. Neurosci.* *17*, 455–462 (2014).
46. Sutton, R. S., & Barto, A. G. (1998). *Reinforcement learning: An introduction*. Cambridge: MIT press.
47. Daw, N. D., Niv, Y., & Dayan, P. (2005). Uncertainty-based competition between prefrontal and dorsolateral striatal systems for behavioral control. *Nat. Neurosci.* *8*, 1704.
48. Huys, Q.J.M., Eshel, N., O’Nions, E., Sheridan, L., Dayan, P., and Roiser, J.P. (2012). Bonsai trees in your head: how the Pavlovian system sculpts goal-directed choices by pruning decision trees. *PLoS Comp. Biol.* *8*, e1002410.
49. Eldar, E., Hauser, T.U., Dayan, P., and Dolan, R.J. (2016) Striatal structure and function predict individual biases in learning to avoid pain. *Proc. Natl. Acad. Sci. USA* *113*, 4812–4817.
50. O’Doherty, J. P., Dayan, P., Friston, K., Critchley, H., & Dolan, R. J. (2003). Temporal difference models and reward-related learning in the human brain. *Neuron* *38*, 329-337.
51. Oostenveld, R., Fries, P., Maris, E. & Schoffelen, J. M. (2011) FieldTrip: open source software for advanced analysis of MEG, EEG, and invasive electrophysiological data. *Comput. intel. Neurosci.* *2011*, 156869.
52. Chang, C.C., and Lin, C.J. (2011). LIBSVM: a library for support vector machines. *ACM T. Intel. Syst. Tec.* *2*, 27.
53. Busch, N. & VanRullen, R. Is visual perception like a continuous flow or a series of snapshots. In: Arstila, V. & Lloyd, D. (Eds.) *Subjective time: The philosophy, psychology, and neuroscience of temporality* (MIT Press, 2014)
54. R Core Team (2018). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL: <https://www.R-project.org/>.
55. Carpenter, B., Gelman, A., Hoffman, M.D., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M., Guo, J., Li, P., & Riddell, A. (2017). Stan: A probabilistic programming language. *Journal of Statistical Software* *76*(1).