

# An ensemble code in medial prefrontal cortex links prior events to outcomes during learning

Silvia Maggi<sup>1</sup>, Adrien Peyrache<sup>2</sup>, and Mark D. Humphries<sup>1\*</sup>

1. Faculty of Biology, Medicine, and Health, University of Manchester, Manchester, UK.
2. Montreal Neurological Institute, McGill University, Montreal, Canada.

\* Contact: [mark.humphries@manchester.ac.uk](mailto:mark.humphries@manchester.ac.uk)

## Abstract

The prefrontal cortex is implicated in learning the rules of an environment through trial and error. But it is unclear how such learning is related to the prefrontal cortex's role in short-term memory. Here we asked if the encoding of short-term memory in prefrontal cortex was used by rats learning decision rules in a Y-maze task. We found that neural ensembles in prefrontal cortex selectively recalled the same pattern of activity after reinforcement for a correct decision. This reinforcement-selective recall only reliably occurred immediately before the abrupt behavioural transitions indicating successful learning of the current rule, and faded quickly thereafter. We could simultaneously decode multiple, retrospective task events from the ensemble activity, suggesting the recalled ensemble activity had multiplexed encoding of prior events. Our results suggest that successful trial-and-error learning is dependent on reinforcement tagging the relevant features of the environment to maintain in prefrontal cortex short-term memory.

## Introduction

Learning the statistical regularities of an environment requires trial and error. But how do we know what is relevant in the environment in order to learn its statistics? In other words: how do we know what to remember? It seems likely that medial prefrontal cortex plays a role here (Euston et al., 2012): it is needed for trial and error learning of correct behavioural strategies (Ragozzino et al., 1999; Ragozzino, 2007; Rich and Shapiro, 2007), neuron and ensemble activity represents abstract and context-dependent information related to the current strategies (Jung et al., 1998; Rich and Shapiro, 2009; Hyman et al., 2012), and changes to ensemble activity tightly correlate with shifts in behavioural strategy (Durstewitz et al., 2010; Karlsson et al., 2012; Powell and Redish, 2016). Moreover, medial prefrontal cortex receives a direct projection from the CA1 field of the hippocampus that may allow the integration of spatial information about the environment (Jones and Wilson, 2005; Hoover and Vertes, 2007; Burton et al., 2009; Benchenane et al., 2010; Spellman et al., 2015). But medial prefrontal cortex also plays a role in short-term and working memory for objects, sequences, and other task features (Miller, 2000; Miller and Cohen, 2001; Baeg et al., 2003; Averbek et al., 2006; Averbek and Lee, 2007; Fujisawa et al., 2008; Jun et al., 2010; Machens et al., 2010; Spellman et al., 2015), upon which

32 successful learning of statistical regularities may depend. It is unknown how relevant  
33 information about the statistics of the environment is tagged for memory in the medial  
34 prefrontal cortex.

35 An hypothesis we consider here is that reinforcement tags relevant choices and features  
36 to remember in order to learn the rules of the environment. If so, then the reliable  
37 appearance of reinforcement-driven short-term memory activity in medial prefrontal cortex  
38 would be predicted during successful learning. As medial prefrontal cortex appears to  
39 encode environmental features and task-related behaviour by ensemble activity (Baeg  
40 et al., 2003; Averbeck and Lee, 2007; Baeg et al., 2007; Sul et al., 2010), any short-term  
41 memory for tagged features would likely be revealed by ensemble activity that was similar  
42 across trials. We thus sought to test the hypothesis that medial prefrontal cortex ensembles  
43 represent a short-term memory of task features and choices that are potentially necessary  
44 for learning from reinforcement.

45 To test this hypothesis, we analysed neural and behavioural data from rats learning  
46 new rules on a Y-maze. We took advantage of a task design in which there was a self-paced  
47 return to the start position of the maze immediately after the delivery or absence of rein-  
48 forcement, yet no explicit working memory component to any of the rules. Consequently  
49 we could examine ensemble activity in medial prefrontal cortex during this self-paced re-  
50 turn and ask whether or not a short-term memory encoding of reinforcement-tagged task  
51 features existed in the absence of overt working memory demands.

52 Here we show that medial prefrontal cortex contains an ensemble code that links prior  
53 events to reinforcement. We show that a neural ensemble activity pattern was specifi-  
54 cally recalled after reinforcement and not after errors. This recall only reliably occurred  
55 in sessions with abrupt shifts in behavioural strategy indicating successful learning, and  
56 not during external shifts in reinforcement contingency, or in other task sessions. From  
57 the activity of the recalled ensemble, we could simultaneously decode retrospective task  
58 parameters and choices in a position-dependent manner. Together, these results show that  
59 learning was preceded by reinforcement-triggered activity of an ensemble that retrospec-  
60 tively and multiply encoded task parameters. They provide a link between the roles of  
61 medial prefrontal cortex in working memory and in rule learning, and suggest that rein-  
62 forcement tags prefrontal cortex-based representations of choices and environment features  
63 that are relevant to trial and error learning of statistical regularities in the world.

## 64 Results

65 In order to address whether and how medial prefrontal cortex neural activity encodes short-  
66 term memory during reinforcement learning, we used medial prefrontal cortex population  
67 recording data previously obtained from a maze-based rule-learning task (Peyrache et al.,  
68 2009). Four rats learnt rules for the direction of the rewarded arm in a Y-shaped maze,  
69 comprising a departure arm and two goal arms with light cues placed next to the reward  
70 ports (Figure 1A). Each session was a single day with approximately 30 minutes of training,  
71 and 30 minutes of pre- and post-training sleep. During training, the rat initiated each trial  
72 from the start of the departure arm; the trial ended when the rat arrived at the reward  
73 point in the goal arm. During the following inter-trial interval the rat made a self-paced  
74 return to the start position after consuming the reward, taking on average 70 s ( $67.8 \pm 5.4$   
75 s, mean  $\pm$  SEM) to complete the return trip. Tetrode recordings from medial prefrontal  
76 cortex were obtained from the very first session in which each rat was exposed to the  
77 maze (Figure 1B). Thus, the combination of a self-paced post-decision period – without  
78 experimenter interference – and neural activity recordings from a naive state allowed us to

79 test for medial prefrontal cortex population activity correlating with short-term memory  
80 during rule learning.

81 After achieving stable performance of the current rule, indicated by 10 contiguous  
82 correct choices, the rule was changed, unsignalled, in sequence: go right; go to the cued  
83 arm; go left; go to the uncued arm. Notably, none explicitly required a working memory  
84 component (such as an alternation rule). In the original study (Peyrache et al., 2009), the  
85 session in which initial learning of each rule occurred was identified posthoc as the first  
86 with three consecutive correct choices followed by 80% performance until the end of the  
87 session; the first of the initial three choices was identified as the learning trial. Ten sessions  
88 met these criteria, and are dubbed here the “learning” sessions. We first confirmed that  
89 these ten learning sessions showed an abrupt transition in behavioural performance (Figure  
90 1C), indicating the step-like change in behaviour commonly seen in successful learning of  
91 contingencies (Gallistel et al., 2004; Aziz-Zadeh et al., 2009; Durstewitz et al., 2010). In  
92 total, we examined 50 sessions, comprising 10 learning sessions, 8 rule change sessions,  
93 and 32 other training sessions (labelled “others” throughout).

#### 94 **Reinforcement-driven recall of ensemble activity during learning**

95 We sought to track reinforcement-driven population activity across the inter-trial inter-  
96 vals within each session, in order to identify signatures of short-term memory encoding.  
97 One signature of similar memory encoding between inter-trial intervals would be the con-  
98 sistent presence of one or more ensembles of neurons with correlated activity. To allow  
99 comparisons between intervals, we thus first identified the core population of neurons in  
100 each session by selecting the neurons that were active in every inter-trial interval. The  
101 proportion of recorded neurons retained in the core population was on average  $74 \pm 2$   
102 % (SEM) across sessions (Figure 2 - figure supplement 1, panel A). No clear difference  
103 in the size of this core population were observed between learning and any other session  
104 type (Figure 2 - figure supplement 1, panel A), suggesting that any potential short-term  
105 memory encoding specific to learning was not then simply a change in the proportion of  
106 active neurons. Rather, any effect of reinforcement on subsequent short-term memory  
107 would have to be encoded in the specific pattern of correlations between the activity of  
108 neurons in the core population.

109 We characterised the pattern of correlations for each inter-trial interval by computing  
110 the pairwise similarity between the Gaussian-convolved spike-trains of neurons in the core  
111 population (we use a Gaussian width of  $\sigma = 100$  ms here, as in the example of Figure 1B;  
112 the effects of varying  $\sigma$  are detailed below). To test if there was one or more reinforcement-  
113 driven ensembles of correlated neurons, we then correlated the core population’s similarity  
114 matrix  $S$  between all inter-trial intervals of a session. The resulting Recall matrix  $R$  showed  
115 where similar patterns of ensemble activity were recalled on different inter-trial intervals  
116 (Figure 2A).

117 We found that patterns of ensemble activity were more similar after correct trials than  
118 after error trials (Figure 2A,B). We observed this preferential post-reinforcement recall  
119 of ensemble activity in the majority of sessions (47/50 sessions; 37/50 had  $p < 0.05$  for  
120 a Kolmogorov-Smirnov test between the distributions of recall values after correct and  
121 after error trials). This result would suggest that reward triggered a specific pattern  
122 of correlated activity during the inter-trial interval. However, we were mindful that the  
123 inter-trial intervals following a correct trial were generally much longer than those following  
124 error trials (correct inter-trial intervals:  $79.1 \pm 6.4$  s; error inter-trial intervals:  $48.4 \pm 3.7$   
125 s), because the animal lingered at the reward location (Figure 2 - figure supplement 1).  
126 This difference in duration could systematically bias estimates of firing correlation, simply

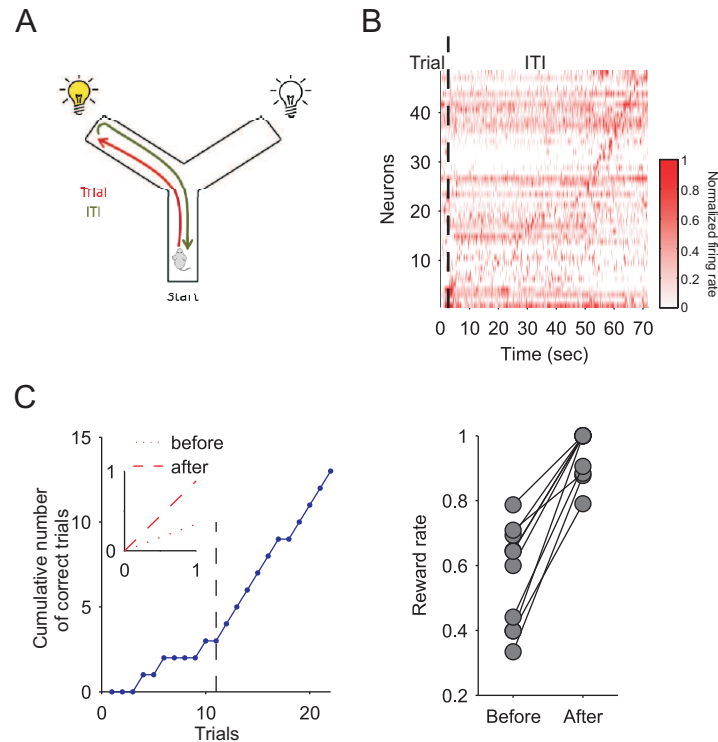


Figure 1: Task and learning sessions.

(A) Schematic representation of the Y maze. The trial starts with the animal at the start of the departure arm, and ends when it reaches the end of the chosen arm. The inter-trial interval (ITI) is a self-paced return back to the start position.

(B) Example medial prefrontal cortex population activity during a trip out and back to the start position. The heatmap shows the spike-trains for all recorded neurons, convolved with a Gaussian of width  $\sigma = 100$  ms. The dashed line separates the trial and inter-trial interval periods. The firing rate of each neuron is a proportion of its peak rate, and neurons are sorted by the time of their peak firing rate.

(C) Learning sessions contain abrupt transitions in performance. Left panel: Learning curve for one example learning session. The cumulative number of correct trials shows a steep increase after the learning trial (black dashed line), indicating the rat had learnt the correct rule. Inset: fitted linear regressions for the cumulative reward before (dotted) and after (dashed) the learning trial, quantifying the large increase in the rate of reward accumulation after the learning trial. Right panel: the rate of reward accumulation before and after the learning trial for every learning session (one pair of symbols per learning session; one session's pair of symbols are obscured). The rate is given by the slopes of the fitted regression lines.

127 because many more spikes would be emitted during post-correct than post-error intervals.  
128 Thus, greater similarity between ensemble activity patterns for post-correct intervals could  
129 simply be due to more reliable estimates of the interval-by-interval correlation matrix. To  
130 control for this, we used shuffled spike-trains to compute the expected matrix of pairwise  
131 similarity due to just the duration of each interval, and from these shuffled-data matrices  
132 we computed the expected recall matrix (Figure 2 - figure supplement 2). Consequently,  
133 by subtracting this expected matrix from the data-derived recall matrix, we obtained a  
134 “residual” recall matrix describing just the similarity between ensemble activity patterns  
135 above those driven by common duration (Figure 2A). We used this residual recall matrix  
136 for all further analyses. With this correction, we still found that patterns of ensemble  
137 activity were more similar after correct trials than after error trials in the majority of  
138 sessions (34/50 sessions; 26/50 had  $p < 0.05$  for a Kolmogorov-Smirnov test between the

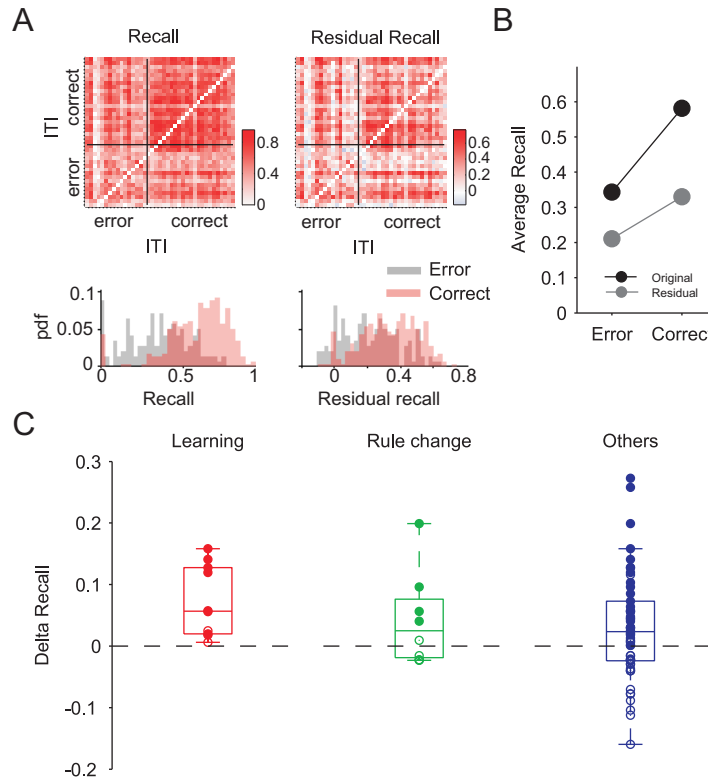


Figure 2: Outcome-selective recall of an ensemble activity pattern is learning-related.

(A) Left: the Recall matrix  $R$  for one example session. The Recall matrix is ordered by the outcome of the trial preceding the inter-trial interval (ITI). Each entry  $R_{ij}$  is the recall value: the similarity between the core population's similarity matrix in intervals  $i$  and  $j$ , a measure of how closely the same ensemble activity pattern was recalled in that pair of intervals. Below we plot the probability density functions for the distribution of recall values, separately for the post-error (red) pairs of intervals (bottom-left block diagonal in the Recall matrix) and for the post-correct (black) pairs of intervals (top right block diagonal in the Recall matrix). Right: the Residual Recall matrix  $R_{resid}$  for the same session, after correction for the effects of interval duration.

(B) The average recall values for post-error and post-correct intervals of the two matrices in panel A. The distribution of recall in the post-correct intervals was higher than in the post-error intervals (K-S test; Recall:  $P < 0.005$ ; Residual recall:  $P < 0.005$ ;  $N(\text{correct}) = 24 \times 24 = 576$ ;  $N(\text{error}) = 17 \times 17 = 238$ .)

(C) The difference in average recall ( $\Delta$  recall) between the post-correct and post-error intervals, sorted by session type. Each dot is one session. Filled circles indicate a positive difference at  $p < 0.05$  between the distributions of recall values in the post-error and post-correct intervals (Kolmogorov-Smirnov test).

139 distributions of residual recall values after correct and after error trials).

140 We then examined how this reinforcement-driven recall of an ensemble activity pattern  
 141 corresponded to the rats' behaviour (Figure 2C). We found that only learning sessions had  
 142 a systematically stronger recall of the same ensemble activity pattern after reinforcement  
 143 (mean  $\Delta$  recall: 0.072). Sessions in which the rule changed did not show a systematic recall  
 144 after reinforcement (mean  $\Delta$  recall = 0.042), ruling out external changes to contingency as  
 145 the driver of the recall effect. Similarly, there was no systematic reinforcement-driven re-  
 146 call in the other sessions (mean  $\Delta$  recall = 0.03), ruling out a general reinforcement-driven  
 147 effect. When we further grouped these other sessions into those with evidence of incremen-  
 148 tal learning and those without, we still did not observe a systematic reinforcement-driven  
 149 recall effect in either group (Figure 2 - figure supplement 2). Finally, we tested the likeli-

150 hood of obtaining ten systematically positive recall sessions by chance if we assumed recall  
151 was randomly distributed across the sessions. We repeatedly chose ten sessions at random  
152 from the 50; repeated 10,000 times, we found a probability of less than 0.003 of randomly  
153 obtaining 10 sessions which each had positive recall. Together, these data show that a sim-  
154 ilar pattern of ensemble activity was only reliably recalled following reinforcement during  
155 the self-driven step-change in behaviour indicative of learning a rule.

156 We asked how the recall of a pattern of ensemble activity was dependent on the tem-  
157 poral precision at which the correlations between neurons were computed. Here, this  
158 precision was determined by the width of the Gaussian convolved with the spike-trains.  
159 We found that the reinforcement-driven recall of an ensemble in learning sessions was  
160 consistent across a wide range of Gaussian widths from 20 ms up to around 140 ms (Fig-  
161 ure 2 - figure supplement 3). Moreover, across the same range of Gaussian widths, we  
162 also consistently found that the recall effect for the learning sessions was greater than for  
163 rule-change or other sessions (Figure 2 - figure supplement 3). The reliable recall down to  
164 20 ms, and the absence of a systematic recall effect for Gaussian widths around 200 ms,  
165 suggests the ensemble was formed by relatively precise correlations between spikes from  
166 different neurons, rather than just rate co-variation.

## 167 **Recall of ensemble activity patterns is specific to retrospective reinforce-** 168 **ment**

169 These results pointed to the hypothesis that, during successful learning of contingency,  
170 the reliable recall of a pattern of ensemble activity is triggered by prior reinforcement.  
171 To test this hypothesis, we asked whether the recalled ensemble was specifically triggered  
172 by reinforcement, and whether it was specific to retrospective rather than prospective  
173 reinforcement.

174 To test if the recall was specifically triggered by reinforcement, we reorganised the  
175 residual recall matrix of each session by either the chosen direction (left/right) or the cue  
176 position (left/right) on the previous trial. We found there was no systematic recall of  
177 ensemble activity patterns evoked by one direction over the other for either the chosen  
178 direction or the cue position (Figure 3A,B). The systematic recall effect during learning  
179 thus appeared to be specific to reinforcement.

180 Modulation of medial prefrontal cortex activity by expected outcome or anticipation  
181 of reinforcement has been repeatedly observed (Daw et al., 2006; Fellows, 2007; Sul et al.,  
182 2010; Kaplan et al., 2017), suggesting the recalled ensemble pattern could instead be a  
183 representation of the expected outcome on the next trial. To test if the recall effect was  
184 specific to retrospective reinforcement, we reordered the residual recall matrices accord-  
185 ing to the reinforcement received in the trial after the inter-trial interval. We found no  
186 systematic recall of an ensemble activity pattern preceding correct trials in any session  
187 type (Figure 3C). In particular, for the learning sessions the systematic recall we observed  
188 for retrospective outcomes was not observed for prospective outcomes (compare Figure  
189 2C), and the magnitude of recall was larger for retrospective than prospective outcomes  
190 across all tested temporal precisions of correlation between spike-trains (Figure 2 - figure  
191 supplement 3, panel D).

192 We were surprised that we could observe such a consistent difference between the ret-  
193 rospective and prospective recall in the learning sessions. By their nature, the learning  
194 sessions tend to be split into a sequence of error trials followed by a sequence of correct  
195 trials (cf Figure 1C), so each trial outcome is frequently preceded and followed by the same  
196 type of outcome. Consequently, whether we split intervals into groups by their following  
197 correct trials or by their preceding correct trials we create similar groups of intervals (and

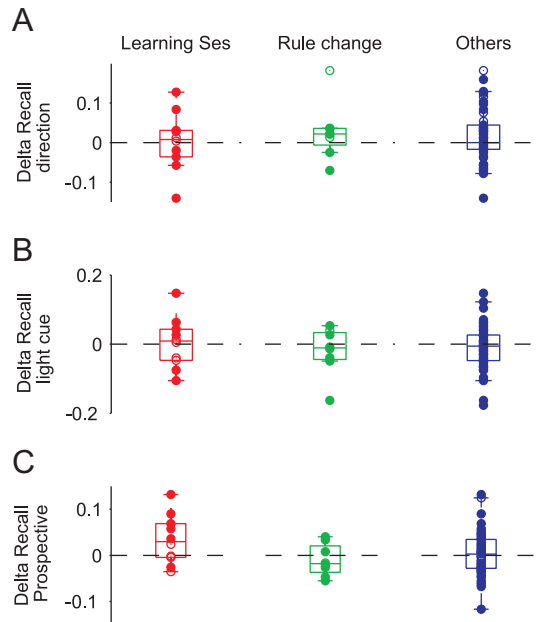


Figure 3: Recalled ensemble activity patterns are outcome-specific and encode retrospective outcome not future choice.

(A) The difference in average recall between intervals after choosing the left respect to the right arm, sorted by session type. Filled circles here and in other panels indicate a significant difference between the distributions of recall values in the two sets of intervals (Kolmogorov-Smirnov test,  $p < 0.05$ ).

(B) As for panel A, but comparing intervals after the light cue appeared at the end of the left or right arm.

(C) The difference in average recall between intervals before error or correct trials, testing for prospective encoding of upcoming choice.

198 similarly for splitting based on error trials). Nonetheless, the systematically stronger ret-  
199 rospective recall across a wide range of timescales, despite the few error trials interspersed  
200 with correct trials, suggests that the recall of ensemble activity is dependent on prior, not  
201 future, reinforcement. (And as we show below, this conclusion is consistent with the com-  
202 plete absence of prospective coding of task elements by the ensemble's activity). Together,  
203 these results support the hypothesis that a specific pattern of ensemble activity triggered  
204 by just-received reinforcement appeared during successful learning of contingency.

## 205 Appearance of the recalled ensemble activity anticipates the behavioural 206 transition

207 This leaves opens the question of whether the appearance of this recalled ensemble pat-  
208 tern is a pre-condition of successful learning, or a read-out of already learnt information.  
209 If a pre-condition, then the recalled ensemble pattern should have appeared before the  
210 transition in behaviour indicating rule acquisition.

211 We thus sought to identify when the recalled ensemble activity pattern first appeared  
212 in each learning session. To do so, we put the recall matrix of each learning session in  
213 trial order (Figure 4A). For each inter-trial interval, we then compared the strength of  
214 recall in the inter-trial intervals before and after that interval (Figure 4B). We used the  
215 inter-trial interval corresponding to the largest difference in recall to identify when the  
216 ensemble activity pattern appeared, as this indicated a step-increase in the similarity of

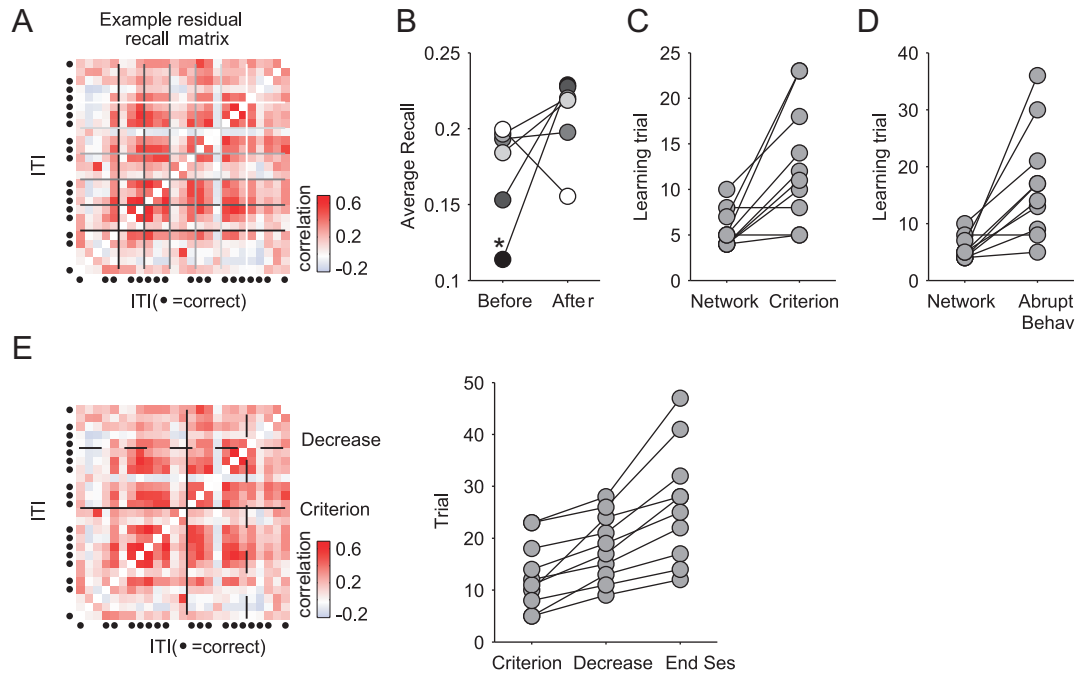


Figure 4: The recalled ensemble activity pattern anticipates behavioural learning.

(A) A residual recall matrix in its temporal order for one example learning session. Columns are ordered from left to right as the first to last inter-trial interval (rows ordered bottom to top). For each inter-trial interval, the distributions of the recall values before and after the selected inter-trial interval were compared (Kolmogorov-Smirnov statistic: see Methods and materials). Each greyscale line corresponds to a selected dividing inter-trial interval plotted in panel B.

(B) For each greyscale line in panel A, the corresponding average recall value before and after the dividing inter-trial interval. The asterisk indicates the inter-trial interval with the largest increase in recall after it, signalling the abrupt appearance of the recalled ensemble pattern.

(C) Comparison of the learning trial identified by the original behavioural criterion and by the abrupt appearance of the recalled ensemble ('Network').

(D) As panel C, but with the behavioural learning trial identified as the trial with the steepest change in the cumulative reward (see Materials and methods).

(E) Testing for decay of the ensemble activity pattern. Left panel: example residual recall matrix in trial order for one learning session. The black solid line is the learning trial, while the dashed line is the identified offset of the recalled ensemble activity pattern. Right panel: For each learning session the learning trial (original criterion) is compared to the identified offset of the ensemble recall, and to the last trial of the session.

217 activity patterns between inter-trial intervals.

218 We found that the recalled ensemble pattern appeared before or approximately si-  
 219 multaneous with the behavioural transition in all sessions (Figure 4C,D). This was true  
 220 whether we used the original behavioural criterion from Peyrache et al. (2009), or our  
 221 more stringent definition of “abrupt” change in the cumulative reward curve (the trial  
 222 corresponding to the greatest change in slope of the reward accumulation curve; see Meth-  
 223 ods). The timing of the appearance of the recalled ensemble pattern was thus consistent  
 224 with it being necessary for successful rule learning.

225 As the change to the ensemble activity was often abrupt and so close to the behavioural  
 226 change, this raised the question of what change to the underlying neural circuit drove this  
 227 change in activity. One possibility would be a physical alteration of connectivity, forming a  
 228 true “structural” cell assembly (Harris, 2005). Alternatively, it could be a temporary effect,  
 229 as might arise from a sustained change in neuromodulation (Durstewitz and Seamans,



230 2002; Benchenane et al., 2011), forming a transient “functional” cell assembly.

231 To decide between these alternatives, we tested for the presence of a long-lasting physi-  
232 cal change by assessing the longevity of the recalled ensemble activity pattern. Specifically,  
233 we tested whether the recall of the ensemble was sustained until the end of the learning  
234 session by performing the onset analysis in reverse (Figure 4E): for each inter-trial interval,  
235 we checked whether the recall after that interval was significantly smaller than before it  
236 (Kolmogorov-Smirnov test; see Materials and methods). We indeed found a statistically  
237 robust fall in the recall of the ensemble activity pattern in every learning session. A strict  
238 ordering was always present: the decay of the recalled ensemble was after the identified  
239 onset of recall, but before the end of the session (Figure 4E), even though we did not  
240 constrain our analysis to this ordering. For the original set of identified learning trials,  
241 the decay trial was always after the learning trial (Figure 4E). (If we used our alternative  
242 learning-trial definition – the trial with the greatest change in reward accumulation – then  
243 7 of the 10 sessions had decay after the learning trial, with 3 sessions showing decay be-  
244 fore it). This analysis indicates the recalled ensemble activity pattern formed transiently  
245 during learning, and decayed quickly after learning was established.

## 246 **Medial prefrontal cortex ensembles had mixed, position-dependent, and** 247 **retrospective encoding of task information**

248 What did the recalled activity pattern encode? Its transient appearance, immediately  
249 before behavioural change but fading before the end of a session, suggests a temporary  
250 representation, akin to short-term memory. That the recalled pattern was triggered only  
251 by prior reinforcement suggests the hypothesis that the recalled ensemble was a working  
252 memory encoding of task features that were potentially relevant for learning. If it was a  
253 working memory for task features, then we should be able to decode prior task information  
254 from ensemble activity.

255 To address this, we assessed our ability to decode prior outcome, choice of direction,  
256 and light cue position from the core population’s activity. As prefrontal cortex activity  
257 encoding often shows broad position dependence (Baeg et al., 2003; Hok et al., 2005;  
258 Spellman et al., 2015), we divided the linearised maze into five equally-spaced sections  
259 (Figure 5A), and represented the core population’s activity in each as the vector of its  
260 neurons’ firing rates in that section. We used these firing rate vectors as inputs to a  
261 cross-validated linear decoder (Figure 5B), and compared their predictive performance to  
262 shuffled data (Materials and methods).

263 We could decode prior outcome, choice of direction and cue position well above chance  
264 performance, and often in multiple contiguous maze positions. We plot the absolute  
265 decoding performance for the “other” sessions in Figure 5C to illustrate that decoding  
266 at some maze positions was near-perfect, with some sessions decoded at 100% accuracy.  
267 The learning and rule-change sessions also had maze positions with near-perfect decoding  
268 across all sessions (Figure 5 - figure supplement 1). Population activity in medial prefrontal  
269 cortex thus robustly encoded multiple task events from the previous trial.

270 We then compared decoding performance between session and rule types. Chance  
271 decoding performance differed between task features (as the randomised light-cue was  
272 counter-balanced across trials, but each rat’s choice and hence outcomes were not), and  
273 between session types and rule types (as rat performance differed between them). Thus  
274 we normalised each decoder’s performance to its own control, and compared this relative  
275 decoding accuracy across sessions and rules (Figure 5D).

276 These comparisons revealed we could decode the prior choice of direction (left or right)  
277 in all types of session and regardless of whether the rule was direction- or cue-based (Figure

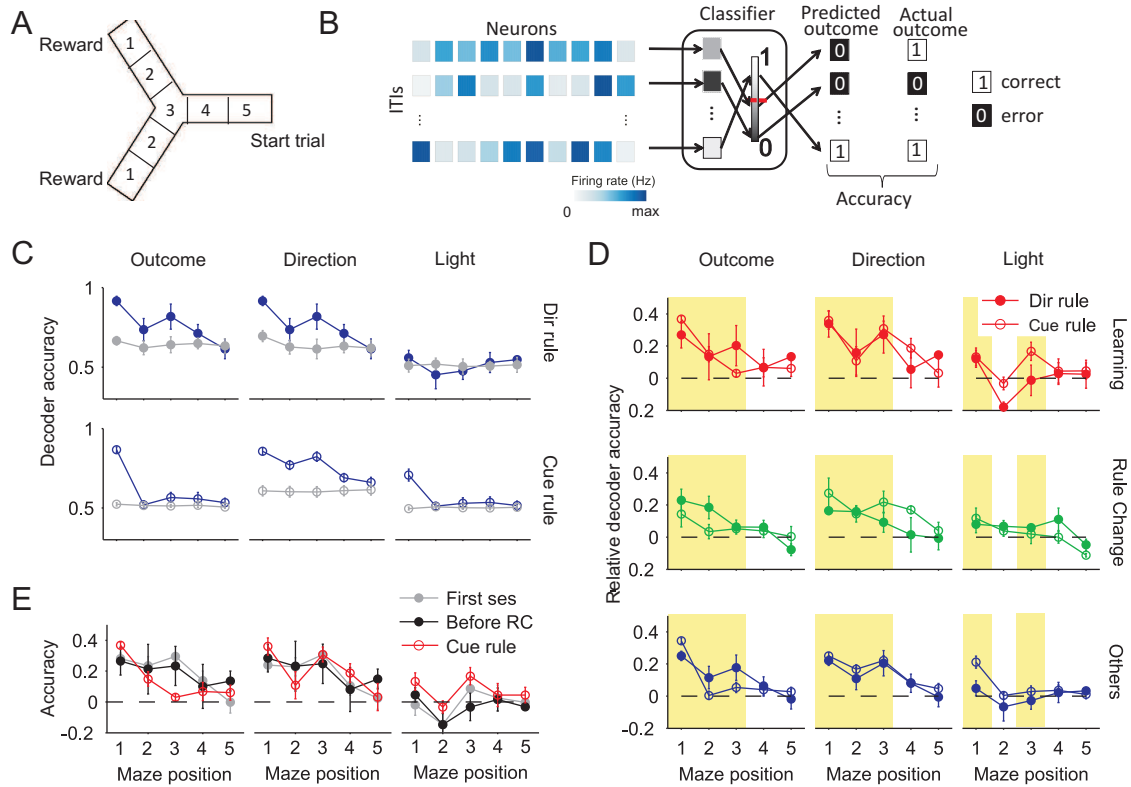


Figure 5: Position-dependent encoding of recent task relevant information.

(A) Graphical representation of the five equally size section of the maze. Position 1 is the goal arm end, with the reward delivery port. Position 3 is the choice point during the trial.

(B) Schematic of decoding task events from the core population's firing rate vector. For each inter-trial interval, and for each of five positions in the maze, the population's firing rate vector is given as input to a linear decoder. The decoder attempts to classify the population vector as belonging to one of two possible labels (error or correct, for prior outcome; left or right, for prior direction choice and prior cue position), given a threshold (red dashed line) on the decoder's output. The accuracy of the decoder is given as the proportion of correctly predicted labels. For robustness, we use cross-validation to create the predictive model for each inter-trial interval; and we compare predictive performance to that of a randomised control: cross-validated classifiers applied to data with permuted labels.

(C) Example decoder accuracy as a function of maze position for the "other" sessions. For these sessions, we show here the absolute decoder performance for each of the three classified features (prior outcome, prior direction choice, and prior light positions), separated by the rule type (direction or light cue-based rules). Each data point is the mean  $\pm$  SEM accuracy at that maze position. Chance levels of performance are plotted in grey, and were defined separately for each session type and each rule type (see Materials and methods). We plot here and in panel D the results for a logistic regression decoder; other decoders are plotted in Figure 5 - figure supplement 2.

(D) Relative decoder accuracy over all session types and rule types. Each data point is the mean  $\pm$  SEM accuracy in excess of chance (0; dashed line) over the indicated combination of session (learning, rule-change, other) and decoded feature. Each panel separately plots the decoding for cued-rules (open symbols) and direction-rules (filled symbols). Highlighted groups of positions indicate consistent departures from chance performance in at least one session type. We replot these results grouped by rule-type in Figure 5 - figure supplement 1.

(E) Decoding of task features at the start of learning. Similar to panel D, each data point is the mean  $\pm$  SEM accuracy in excess of chance. In each panel, the grey line gives the accuracy over the first session of each animal, the black line the accuracy over all the sessions before the first rule change, and the red line gives the accuracy over the first light cue session for each animal.

278 5D). Decoding of direction choice was robustly above chance while the rats moved from  
279 the end of the goal arm back to the maze's choice point (highlighted yellow); on cued-  
280 rule sessions, this decoding extended almost all the way back to the start position of the  
281 departure arm. Accurate decoding of direction choice could be observed from the very first  
282 session of each rat, and consistently across sessions before the first rule change (Figure  
283 5E). These results indicated that medial prefrontal cortex always maintained a memory  
284 of prior choice, and did not need to learn to encode this task feature.

285 Similarly, we could decode the prior outcome (correct or error) in all types of session  
286 and regardless of whether the rule was direction- or cue-based (Figure 5D). Decoding of  
287 outcome was notably stronger at the end of the goal arm, where the reward was delivered,  
288 but could also be decoded above chance while the rats traversed the maze back to the start  
289 position (highlighted yellow). Nonetheless, decoding of outcome was again present from  
290 the very first session (Figure 5E). These results indicated that medial prefrontal cortex  
291 always encoded the trial's outcome, and did not need to learn to encode this task feature.

292 By contrast to the encoding of prior direction and outcome, we could only reliably  
293 decode the prior cue position in two specific locations (Figure 5D). The prior cue position  
294 was consistently encoded at the end of the goal arm for both cue and direction rules,  
295 likely corresponding to whether or not the light was on at the rat's position. But the  
296 only sustained encoding of prior cue position while the rat traversed the maze was during  
297 learning sessions for cue-based rules (yellow highlighted position and red open circles in  
298 Figure 5D). There was no sustained encoding during learning sessions of direction rules  
299 (red filled circles in Figure 5D). And this sustained encoding did not appear in the first  
300 session, nor in any session before the first change to the cue-based rule (Figure 5E).  
301 Consequently, these results suggest that only in learning sessions did the core population  
302 encode the memory of the prior cue position, and only when relevant to the learnt rule.

303 Strikingly, we found that decoding of prospective choice or outcome on the next trial  
304 was at chance levels throughout the inter-trial interval (Figure 6). These results were  
305 consistent with our finding that the ensemble activity pattern preceding correct trials  
306 was not systematically recalled (Figure 3C). They also show that the decoding of prior  
307 task features from the core populations' activity was non-trivial. The only above-chance  
308 decoding of prospective information was observed for direction-based rules, where we found  
309 that decoding of future choice and outcome was above chance level only for learning  
310 sessions and only at position 5, where the animal make a U-turn before starting the new  
311 trial. This suggests that medial prefrontal cortex activity around the start of the trial  
312 could also be related to the upcoming decision when the task rule is successfully learnt;  
313 future work will explore this idea.

314 Collectively, the decoder analysis showed we could decode multiple task features from  
315 the immediate past from population activity in the medial prefrontal cortex – in some  
316 cases, perfectly (Figure 5C, and Figure 5 - figure supplement 1) – but not the immediate  
317 future. Moreover, our ability to decode these prior task features was consistent across a  
318 range of tested decoders, as was the sustained encoding of prior cue position only during  
319 the learning of cue-based rules (Figure 5 - figure supplement 2). Thus, we suggest that the  
320 specific pattern of recalled ensemble activity triggered by reinforcement is the repeated  
321 synchronisation of the multiplexed encoding of prior choice, outcome, and cue position  
322 relevant to learning the current rule.

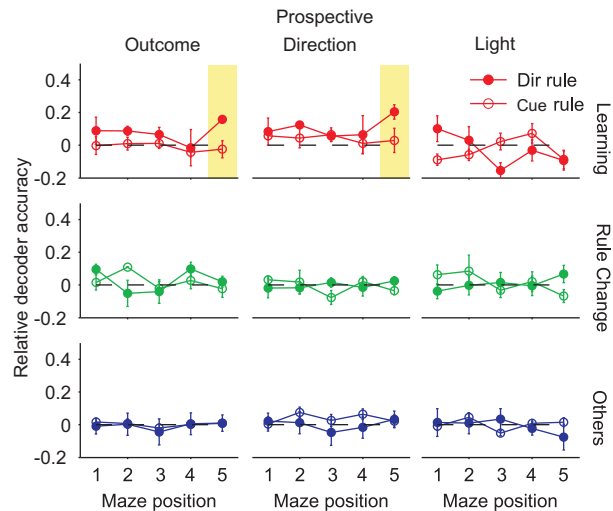


Figure 6: Prospective encoding of outcome and choice for predictable rule. We plot here the relative decoder accuracy over all session types and rule types for features on the immediately following trial. Compare to Figure 5D. Each data point is the mean  $\pm$  SEM accuracy in excess of chance (dashed line) over the indicated combination of session (learning, rule-change, other) and decoded feature. Each panel separately plots the decoding for cued-rules (open symbols) and direction-rules (filled symbols). As a sanity check that our cross-validation of the decoder and shuffled controls were working, we also decoded the prospective light position: as this was randomised, the ensemble activity could not predict its position and so should only have been decoded at chance levels – which it was. We replot these results grouped by rule-type in Figure 5 - figure supplement 1.

## 323 Discussion

324 We sought to understand how short-term memory in medial prefrontal cortex may support  
 325 the trial-and-error learning of rules from a naive state. To do so, we analysed population  
 326 activity in medial prefrontal cortex from rats learning rules on a maze, and asked if the  
 327 activity during the inter-trial interval carried signatures of short-term memory for rule-  
 328 relevant features of the task. We found that a specific pattern of ensemble activity was  
 329 recalled only after reinforced trials, and only reliably during sessions in which the rats  
 330 learnt the current rule for the first time. This dependence on prior outcome, and the  
 331 transient appearance of the ensemble activity pattern, was consistent with a short-term  
 332 memory encoding, rather than a persistent change to the underlying neural circuit.

333 We could robustly decode prior outcome and direction choice from ensemble activity  
 334 across all sessions, but found that encoding of the prior cue position was specific to learning  
 335 sessions for the cue-based rules. This suggests that the recalled ensemble is a repeated  
 336 synchronisation of multiple encodings across the neural population, with rule-appropriate  
 337 suppression or enhancement of cue encoding. We thus propose that reinforcement tags  
 338 features to sustain in medial prefrontal cortex working memory, and does this by reliably  
 339 triggering a specific pattern of ensemble activity that jointly encodes relevant task features.

## 340 Ensemble recall precedes behavioural learning

341 We only reliably observed the recall of the same pattern of correlation between neurons  
 342 following correct trials during learning. This pattern of correlation reliably appeared before  
 343 or simultaneously with the step-like change in behaviour indicating rule acquisition. The  
 344 timing thus suggests a causal link between the appearance of the recalled ensemble activity  
 345 pattern and successful learning of the correct strategy for the current rule.

346 Our results support recent studies of prefrontal cortex population activity that re-  
347 ported how the pattern of population activity in rodent prefrontal cortex changes with or  
348 immediately prior to an internally-driven shift in behavioural strategy (Durstewitz et al.,  
349 2010; Powell and Redish, 2016). We extend these prior results in three ways. First, prior  
350 work has studied the scenarios where animals well-trained on one contingency experienced  
351 a change in that contingency. Here, we have shown that such abrupt shifts in population  
352 activity patterns can occur from the naive state. Consequently, they encode initial acqui-  
353 sition as well as uncertainty (Karlsson et al., 2012). Second we have shown that such an  
354 abrupt shift in population activity happens for a putative working memory representa-  
355 tion. Third, we have shown that this shift is selectively triggered by prior reinforcement.  
356 Nonetheless, our results add to the growing evidence that an abrupt shift in prefrontal  
357 cortex population activity is a necessary condition for the successful acquisition of a new  
358 behavioural strategy.

### 359 **Functional cell assemblies are potentially necessary for learning but not** 360 **performance**

361 That the recalled ensembles only appeared around clear episodes of behavioural learning  
362 means they are thus candidate cell assemblies (Harris, 2005): an ensemble that appeared  
363 during the course of learning. We distinguished here between structural and functional  
364 cell assemblies. In a structural assembly, the ensemble's activity pattern is formed by  
365 some underlying physical change, such as synaptic plasticity of the connections between  
366 and into the neurons of the ensemble (Harris, 2005; Holtmaat and Caroni, 2016), and is  
367 thus a permanent change. In a functional assembly, the ensemble's activity pattern is  
368 formed by some temporary modulation of existing connections - e.g. by new input or  
369 neuromodulation (Benchenane et al., 2011), and is thus a temporary change. Our analysis  
370 suggested that the recalled ensembles were a functional assembly, as they decayed before  
371 the end of the session in which they appeared, often decaying soon after the learning  
372 trial itself. We thus propose that this short-term memory ensemble is necessary only for  
373 the successful trial-and-error learning of a new rule, and not for the ongoing successful  
374 performance of that rule.

### 375 **Encoding in prefrontal cortex from the naive state**

376 Consistent with prior reports of mixed selectivity in prefrontal cortex (Jung et al., 1998;  
377 Jun et al., 2010; Rigotti et al., 2013), we could decode multiple task features from the  
378 joint activity of a small population of neurons. Extending these reports, we showed here  
379 that these encodings were position dependent, and that this encoding was exclusively ret-  
380 rospective during the inter-trial interval - despite there being no explicit working memory  
381 component to the rules. Our data thus show a short-term memory for multiplexed task  
382 features even in the absence of overt working memory demands.

383 One of our more unexpected findings was that we could reliably decode both the  
384 prior choice of direction and the prior trial's outcome across all sessions, regardless of  
385 whether they contained clear learning, externally-imposed rule changes, or neither these  
386 events. Our decoder used the vector of firing rates at a given maze position as input.  
387 Consequently, our ability to decode binary labels of prior events (correct/error trials or  
388 left/right locations) implies that there were well separated firing rate vectors for each of  
389 these labels. But this does not mean the neurons' firing rates were consistently related  
390 for a given label (such as a prior choice of the left arm of the maze). Indeed, it could  
391 imply anything from the two labels being encoded by the only two vectors of firing rates

392 that ever appeared, to the two labels being encoded by two distinct groups of neurons  
393 whose firing rates within each group were never correlated. The reliable appearance of the  
394 same pattern of pairwise correlations only during learning thus implies that only during  
395 these sessions was the firing rate vector reliably correlated. This suggests that learning to  
396 synchronise the encoded features, and not the learning of the encoding itself, is necessary  
397 for acquiring of a new rule.

398 An interesting detail with potentially broad implications is that we could decode both  
399 the choice of prior direction and prior outcome from the very first session that each rat  
400 experienced the Y-maze. Either this implies that medial prefrontal cortex learnt rep-  
401 resentations of direction and outcome so fast that they were able to make a significant  
402 contribution to decoding by population activity within the very first session. Or it im-  
403 plies that medial prefrontal cortex does not need to learn representations of direction and  
404 outcome, meaning that such encoding is always present. Future work is needed to distin-  
405 guish which of the broad spectrum of features encoded by the prefrontal cortex are either  
406 consistently present or learnt according to task demands. Demarcating the classes of fea-  
407 tures that the prefrontal cortex innately or learns to remember would further advance our  
408 understanding of its contribution to adaptive behaviour.

## 409 **Acknowledgments**

410 We thank Matt Jones for comments on the manuscript, and the Humphries' lab (Abhinav  
411 Singh, Javier Caballero, Mat Evans) for discussion. M.D.H and S.M. were supported  
412 by a Medical Research Council Senior non-Clinical Fellowship award (MR/J008648/1) to  
413 M.D.H.

## 414 **Materials and methods**

### 415 **Task description and electrophysiological data**

416 For full details on training, spike-sorting, and histology see (Peyrache et al., 2009). Four  
417 Long-Evans male rats with implanted tetrodes in prelimbic cortex were trained on a Y-  
418 maze task (Figure 1A). Each recording session consisted of a 20-30 minute sleep or rest  
419 epoch, in which the rat remained undisturbed in a padded flowerpot placed on the central  
420 platform of the maze, followed by a training epoch, in which the rat performed for 20-40  
421 minutes, and then by a second 20-30 minute sleep or rest epoch. Periods of slow-wave  
422 sleep were detected offline automatically from local field potential recordings (details in  
423 Peyrache et al., 2009).

424 During training, every trial started when the rat left the beginning of the start arm  
425 and finished when the rat reached the end of one of the choice arms. A correct choice of  
426 arm was rewarded with drops of flavoured milk. Each inter-trial interval lasted from the  
427 end-point of the trial until the rat made its self-paced return to the beginning of the start  
428 arm.

429 Each rat had to learn the current rule by trial-and-error. The rules were sequenced to  
430 ensure cross-modal shifts: go to the right arm; go to the cued arm; go to the left arm; go  
431 to the uncued arm. To maintain consistent context across all sessions, the light cues were  
432 lit in a pseudo-random sequence across trials, whether they were relevant to the rule or  
433 not.

434 The data analysed here were from a total set of 53 experimental sessions taken from  
435 the study of Peyrache et al. (2009), representing a set of training sessions from naive until

436 either the final training session, or until choice became habitual across multiple consecutive  
437 sessions (consistent selection of one arm that was not the correct arm). In this data-set,  
438 each rat learnt at least two rules, and the four rats respectively contributed 14, 14, 11, and  
439 14 sessions. We used 50 sessions here, omitting one session for missing position data, one in  
440 which the rat always chose the right arm (in a dark arm rule) preventing further decoding  
441 analyses (see below), and one for missing spike data in a few trials. Tetrode recordings  
442 were obtained from the first session for each rat. They were spike-sorted only within each  
443 recording session for conservative identification of stable single units. In the sessions we  
444 analyse here, the populations ranged in size from 15-55 units. Spikes were recorded with  
445 a resolution of 0.1 ms. Simultaneous tracking of the rat's position was recorded at 30 Hz.

446 In order to identify ensembles and track them over each session, we first selected the  
447  $N$  neurons that were active in all the inter-trial intervals. The  $N$  spike-trains of this  
448 core population were convolved with a Gaussian ( $\sigma = 100$  ms) to obtain a spike-density  
449 function  $f_k$  for the  $k$ th spike-train. All the recall analysis was repeated for different  
450 Gaussian widths ranging from 20 ms to 240 ms (Figure 2 - figure supplement 3). Each  
451 spike-train was then Z-scored to obtain a normalised spike-density function  $f_k^*$  of unit  
452 variance:  $f_k^* = (f_k - \langle f_k \rangle) / \sigma_k$ , where  $\langle f_k \rangle$  is the mean of  $f_k$ , and  $\sigma_k$  its standard deviation.

### 453 Testing for reinforcement-driven ensembles

454 To compare the core population's pattern of activity across the session, for each inter-  
455 trial interval  $i$  we first computed a pairwise similarity matrix  $S_i$  between the spike-density  
456 functions for all  $N$  neurons. Similarity here was the rectified correlation coefficient, re-  
457 taining all positive values, and setting all negative values to zero. We did this because,  
458 as detailed below, we needed to decompose the pairwise measurements into two additive  
459 contributions: we thus restricted pairwise measurements to the positive regime so that the  
460 difference in contributions lay on the interval  $[-1,1]$ , and so that two negative contributions  
461 could not sum to a positive contribution.

462 We then compared the core population's correlation patterns between inter-trial inter-  
463 vals  $i$  and  $j$  by computing the pairwise similarity between  $S_i$  and  $S_j$ . By comparing all  
464 pairs of inter-trial intervals, we thus formed the Recall matrix  $R$ , capturing the similarity  
465 of activity patterns between all inter-trial intervals.

466 We grouped the entries of  $R$  into two groups according to the same type of inter-  
467 trial interval - predominantly whether they were intervals following correct or following  
468 error trials. These created the block diagonals  $R_1$  and  $R_2$  (such as  $R_{error}$  and  $R_{correct}$ ,  
469 as illustrated in Figure 2A). We summarised the recall between groups by computing the  
470 mean of each block. We detected statistically meaningful differences by computing the  
471 Kolmogorov-Smirnov test for a difference between the distributions of values in the two  
472 blocks.

473 In the main text, we report that there is higher average similarity in  $R_{correct}$  than  
474  $R_{error}$  in many sessions. However, there was a strong tendency for inter-trial intervals  
475 following correct trials to be longer in duration than inter-trial intervals following error  
476 trials (Figure 2 - figure supplement 1), and so the estimates of pairwise similarity may  
477 be biased. In order to dissect the contribution of the different durations we defined a  
478 null model. For each session we defined a shuffled Recall matrix  $\hat{R}$  obtained from the  
479 average of 1000 Recall matrices computed on shuffled spike trains, with each shuffled  
480 spike-train keeping its inter-spike interval distribution fixed. In this way we destroyed  
481 any task-specific temporal pattern of the spike train and we quantify the contribution to  
482 pairwise similarity due solely to the length of the inter-trial interval. Our final residual  
483 Recall matrix  $\tilde{R} = R - \hat{R}$  is obtained as the difference between the Recall matrix and the

484 average shuffled Recall matrix (Figure 2A; Figure 2 - figure supplement 1).

485 For the Residual Recall matrix, we summarised and tested the differences between the  
486 two groups (such as post-error and post-correct inter-trial intervals) in the same way as  
487 detailed above, given the new block diagonals  $\tilde{R}_1$  and  $\tilde{R}_2$ . When grouping by session type,  
488 we plotted the difference between the block diagonals' means as Delta Recall =  $\text{mean}(\tilde{R}_1)$   
489 -  $\text{mean}(\tilde{R}_2)$ .

## 490 Behavioural analysis

491 A learning trial was defined following the criteria of the original study (Peyrache et al.,  
492 2009) as the first of three correct trials after which the performance was at least 80%  
493 correct for the remainder of a session. Only ten sessions contained a trial which met  
494 these criteria, and so were labelled “learning” sessions. We checked that these identified  
495 trials corresponded to an abrupt change in behaviour by computing the cumulative reward  
496 curve, then fitting a piecewise linear regression model: a robust regression line fitted to  
497 the curve before the learning trial, and another fitted to the curve after the learning trial.  
498 The slopes of the two lines thus gave us the rate of reward accumulation before ( $r_{before}$ )  
499 and after ( $r_{after}$ ) the learning trial.

500 To identify other possible learning trials within the learning session, we fitted this  
501 piecewise linear regression model to each trial in turn (allowing a minimum of 5 trials  
502 before and after each tested trial). We then found the trial at which the increase in  
503 slope ( $r_{after} - r_{before}$ ) was maximised, indicating the point of steepest inflection in the  
504 cumulative reward curve. The two sets of learning trials largely agreed: we checked our  
505 results using this set too.

506 Amongst the other sessions, we searched for signs of incremental learning by again  
507 fitting the piecewise linear regression model to each trial in turn, and looking for any trial  
508 for which ( $r_{after} - r_{before}$ ) was positive. We found 22 sessions falling in this category in  
509 addition to the 10 learning sessions. We called those new sessions Minor-learning (Figure  
510 2 - Supplement figure 2).

## 511 Testing the onset and offset of recall

512 In order to identify when the recalled ensemble activity pattern first appeared in a learning  
513 session, we arranged its Residual Recall matrix in trial order. For each trial in turn (with a  
514 minimum of 3 trials before and 5 after), we formed the block diagonals  $R_{before}$  and  $R_{after}$   
515 (see Figure 4A), respectively giving all pairwise recall scores between inter-trial intervals  
516 before and after that trial. The distance between recall before and after was measured  
517 using the Kolmogorov-Smirnov statistic: the maximum distance between the empirical  
518 cumulative distributions of  $R_{before}$  and  $R_{after}$ . The trial that had the maximum positive  
519 distance (an increase in recall from  $R_{before}$  to  $R_{after}$ ) and had  $P < 0.05$  was identified as  
520 the onset of the recalled activity pattern. Similarly, the trial with the maximum negative  
521 distance that corresponded to a decrease in recall from  $R_{before}$  to  $R_{after}$  and had  $P < 0.05$   
522 was identified as the offset of the recalled activity pattern. In all learning sessions we  
523 observed a strict ordering of onset occurring before offset, and both occurring before the  
524 final tested trial of the session.

## 525 Decoder analysis

526 To test whether it was possible to predict task-relevant information in a position-dependent  
527 manner from the core population's activity we trained and tested a range of linear decoders



528 (Hastie et al., 2009). In the main text we report the results obtained using a logistic  
529 regression classifier, as this is perhaps the easiest classifier to interpret.

530 We first linearised the maze in five equally-sized sections, with the central section  
531 covering the choice point of the maze. During each inter-trial interval, we computed  
532 the  $N$ -length firing rate vector  $R^p$ , whose each element  $r_j^p$  is the firing rate of the  $j$ th  
533 core population neuron at position  $p$ . For each session of  $T$  inter-trial intervals and each  
534 section of the maze  $p$ , the set of population firing rate vectors  $R^p(1), \dots, R^p(T)$  was then  
535 used to train a linear decoder to classify the relevant binary task information, either: the  
536 previous trial's outcome (labels: 0,1), the previously chosen arm (labels: left,right), or the  
537 previous position of the light cue (labels: left,right). (We also trained all decoders on the  
538 next outcome, arm choice, and light position to test for prospective encoding). To avoid  
539 overfitting, we used leave-one-out cross-validation, where each inter-trial interval was held  
540 out in turn as the test target and the decoder was trained on the  $T - 1$  remaining inter-  
541 trial intervals. The accuracy of the decoder for position  $p$  in a given session was thus the  
542 proportion of correctly predicted labels over the  $T$  held out test inter-trial intervals.

543 Because the frequency of outcomes and arm choices were due to the rat's behaviour,  
544 chance proportions of correctly decoding labels was not 50%. To establish chance perfor-  
545 mance for each decoding, we fitted the same cross-validated classifier on the same set of  
546 firing rate vectors at each position, but using shuffled labels across the inter-trial intervals  
547 (for example, we shuffled the outcomes of the previous trial randomly). We repeated the  
548 shuffling and fitting 50 times. For displaying the results in Figure 5, we subtracted the  
549 mean of the shuffled results from the true decoding performance. Separate results for the  
550 true and shuffled decoders are plotted in Figure 5 - figure supplement 1, panel A.

551 We report in the main text the results of using a logistic regression classifier. To  
552 check the robustness of our results, we also tested three further linear decoders: linear  
553 discriminant analysis; (linear) support vector machines; and a nearest neighbours classifier.  
554 Each of these showed similar decoding performance to the logistic regression classifier  
555 (Figure 5 - figure supplement 2).

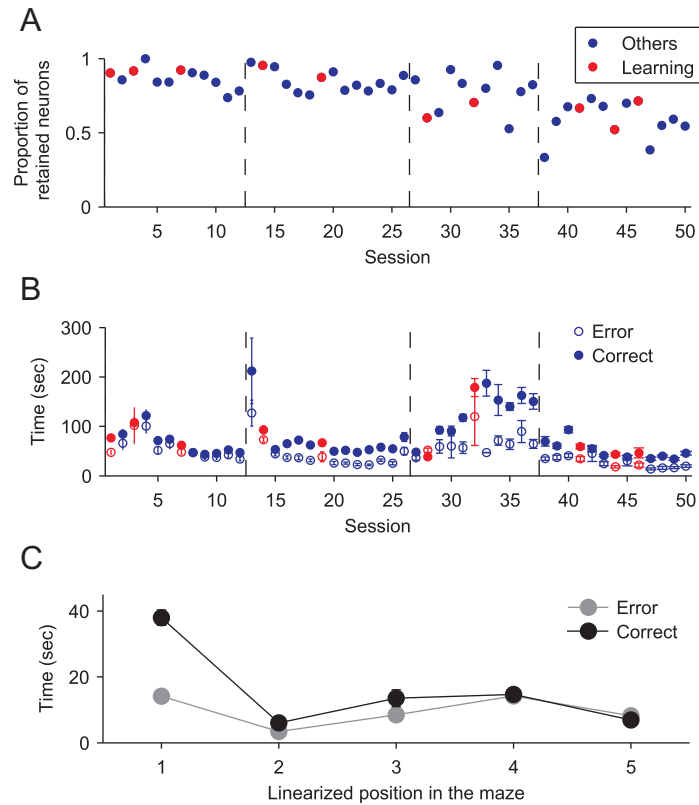
## 556 References

- 557 Averbeck, B. B. and Lee, D. (2007). Prefrontal neural correlates of memory for sequences.  
558 *J Neurosci*, 27:2204–2211.
- 559 Averbeck, B. B., Sohn, J.-W., and Lee, D. (2006). Activity in prefrontal cortex during  
560 dynamic selection of action sequences. *Nat Neurosci*, 9:276–282.
- 561 Aziz-Zadeh, L., Kaplan, J. T., and Iacoboni, M. (2009). aha!: The neural correlates of  
562 verbal insight solutions. *Human brain mapping*, 30(3):908–916.
- 563 Baeg, E. H., Kim, Y. B., Huh, K., Mook-Jung, I., Kim, H. T., and Jung, M. W. (2003).  
564 Dynamics of population code for working memory in the prefrontal cortex. *Neuron*,  
565 40(1):177–188.
- 566 Baeg, E. H., Kim, Y. B., Kim, J., Ghim, J.-W., Kim, J. J., and Jung, M. W. (2007).  
567 Learning-induced enduring changes in functional connectivity among prefrontal cortical  
568 neurons. *Journal of Neuroscience*, 27(4):909–918.
- 569 Benchenane, K., Peyrache, A., Khamassi, M., Tierney, P. L., Gioanni, Y., Battaglia, F. P.,  
570 and Wiener, S. I. (2010). Coherent theta oscillations and reorganization of spike timing  
571 in the hippocampal- prefrontal network upon learning. *Neuron*, 66(6):921–936.

- 572 Benchenane, K., Tiesinga, P. H., and Battaglia, F. P. (2011). Oscillations in the prefrontal  
573 cortex: a gateway to memory and attention. *Curr Opin Neurobiol*, 21(3):475–485.
- 574 Burton, B. G., Hok, V., Save, E., and Poucet, B. (2009). Lesion of the ventral and  
575 intermediate hippocampus abolishes anticipatory activity in the medial prefrontal cortex  
576 of the rat. *Behav Brain Res*, 199(2):222–234.
- 577 Daw, N. D., O’doherly, J. P., Dayan, P., Seymour, B., and Dolan, R. J. (2006). Cortical  
578 substrates for exploratory decisions in humans. *Nature*, 441(7095):876–879.
- 579 Durstewitz, D. and Seamans, J. K. (2002). The computational role of dopamine D1  
580 receptors in working memory. *Neural Netw*, 15(4-6):561–572.
- 581 Durstewitz, D., Vittoz, N. M., Floresco, S. B., and Seamans, J. K. (2010). Abrupt transi-  
582 tions between prefrontal neural ensemble states accompany behavioral transitions during  
583 rule learning. *Neuron*, 66(3):438–448.
- 584 Euston, D. R., Gruber, A. J., and McNaughton, B. L. (2012). The role of medial prefrontal  
585 cortex in memory and decision making. *Neuron*, 76(6):1057–1070.
- 586 Fellows, L. K. (2007). Advances in understanding ventromedial prefrontal function the  
587 accountant joins the executive. *Neurology*, 68(13):991–995.
- 588 Fujisawa, S., Amarasingham, A., Harrison, M. T., and Buzsaki, G. (2008). Behavior-  
589 dependent short-term assembly dynamics in the medial prefrontal cortex. *Nat Neurosci*,  
590 11(7):823–833.
- 591 Gallistel, C. R., Fairhurst, S., and Balsam, P. (2004). The learning curve: implications of  
592 a quantitative analysis. *Proceedings of the national academy of Sciences of the united  
593 States of america*, 101(36):13124–13131.
- 594 Harris, K. D. (2005). Neural signatures of cell assembly organization. *Nature Reviews  
595 Neuroscience*, 6(5):399–407.
- 596 Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The Elements of Statistical Learning*.  
597 Springer, Berlin.
- 598 Hok, V., Save, E., Lenck-Santini, P. P., and Poucet, B. (2005). Coding for spatial goals  
599 in the prelimbic/infralimbic area of the rat frontal cortex. *Proc Natl Acad Sci U S A*,  
600 102(12):4602–4607.
- 601 Holtmaat, A. and Caroni, P. (2016). Functional and structural underpinnings of neuronal  
602 assembly formation in learning. *Nature Neuroscience*, 19:1553–1562.
- 603 Hoover, W. B. and Vertes, R. P. (2007). Anatomical analysis of afferent projections to the  
604 medial prefrontal cortex in the rat. *Brain Struct Funct*, 212:149–179.
- 605 Hyman, J. M., Ma, L., Balaguer-Ballester, E., Durstewitz, D., and Seamans, J. K. (2012).  
606 Contextual encoding by ensembles of medial prefrontal cortex neurons. *Proceedings of  
607 the National Academy of Sciences*, 109(13):5086–5091.
- 608 Jones, M. W. and Wilson, M. A. (2005). Theta rhythms coordinate hippocampal-prefrontal  
609 interactions in a spatial memory task. *PLoS Biol*, 3(12):e402.

- 610 Jun, J. K., Miller, P., Hernandez, A., Zainos, A., Lemus, L., Brody, C. D., and Romo, R.  
611 (2010). Heterogenous population coding of a short-term memory and decision task. *J*  
612 *Neurosci*, 30(3):916–929.
- 613 Jung, M. W., Qin, Y., McNaughton, B. L., and Barnes, C. A. (1998). Firing characteristics  
614 of deep layer neurons in prefrontal cortex in rats performing spatial working memory  
615 tasks. *Cereb Cortex*, 8(5):437–450.
- 616 Kaplan, R., King, J., Koster, R., Penny, W. D., Burgess, N., and Friston, K. J. (2017).  
617 The neural representation of prospective choice during spatial planning and decisions.  
618 *PLOS Biology*, 15(1):e1002588.
- 619 Karlsson, M. P., Tervo, D. G., and Karpova, A. Y. (2012). Network resets in medial  
620 prefrontal cortex mark the onset of behavioral uncertainty. *Science*, 338(6103):135–139.
- 621 Machens, C. K., Romo, R., and Brody, C. D. (2010). Functional, but not anatomical,  
622 separation of “what” and “when” in prefrontal cortex. *J Neurosci*, 30:350–360.
- 623 Miller, E. K. (2000). The prefrontal cortex and cognitive control. *Nature reviews neuro-*  
624 *science*, 1(1):59–65.
- 625 Miller, E. K. and Cohen, J. D. (2001). An integrative theory of prefrontal cortex function.  
626 *Annual review of neuroscience*, 24(1):167–202.
- 627 Peyrache, A., Khamassi, M., Benchenane, K., Wiener, S. I., and Battaglia, F. P. (2009).  
628 Replay of rule-learning related neural patterns in the prefrontal cortex during sleep.  
629 *Nature Neuroscience*, 12:919–926.
- 630 Powell, N. J. and Redish, A. D. (2016). Representational changes of latent strategies in  
631 rat medial prefrontal cortex precede changes in behaviour. *Nature Communications*, 7.
- 632 Ragozzino, M. E. (2007). The contribution of the medial prefrontal cortex, orbitofrontal  
633 cortex, and dorsomedial striatum to behavioral flexibility. *Annals of the New York*  
634 *Academy of Sciences*, 1121(1):355–375.
- 635 Ragozzino, M. E., Wilcox, C., Raso, M., and Kesner, R. P. (1999). Involvement of rodent  
636 prefrontal cortex subregions in strategy switching. *Behav Neurosci*, 113(1):32–41.
- 637 Rich, E. L. and Shapiro, M. (2009). Rat prefrontal cortical neurons selectively code  
638 strategy switches. *Journal of Neuroscience*, 29(22):7208–7219.
- 639 Rich, E. L. and Shapiro, M. L. (2007). Prelimbic/infralimbic inactivation impairs mem-  
640 ory for multiple task switches, but not flexible selection of familiar tasks. *J Neurosci*,  
641 27(17):4747–4755.
- 642 Rigotti, M., Barak, O., Warden, M. R., Wang, X.-J., Daw, N. D., Miller, E. K., and Fusi,  
643 S. (2013). The importance of mixed selectivity in complex cognitive tasks. *Nature*,  
644 497(7451):585–590.
- 645 Spellman, T., Rigotti, M., Ahmari, S. E., Fusi, S., Gogos, J. A., and Gordon, J. A. (2015).  
646 Hippocampal-prefrontal input supports spatial encoding in working memory. *Nature*,  
647 522(7556):309–314.
- 648 Sul, J. H., Kim, H., Huh, N., Lee, D., and Jung, M. W. (2010). Distinct roles of rodent  
649 orbitofrontal and medial prefrontal cortex in decision making. *Neuron*, 66(3):449–460.

650 **Supplementary figures**



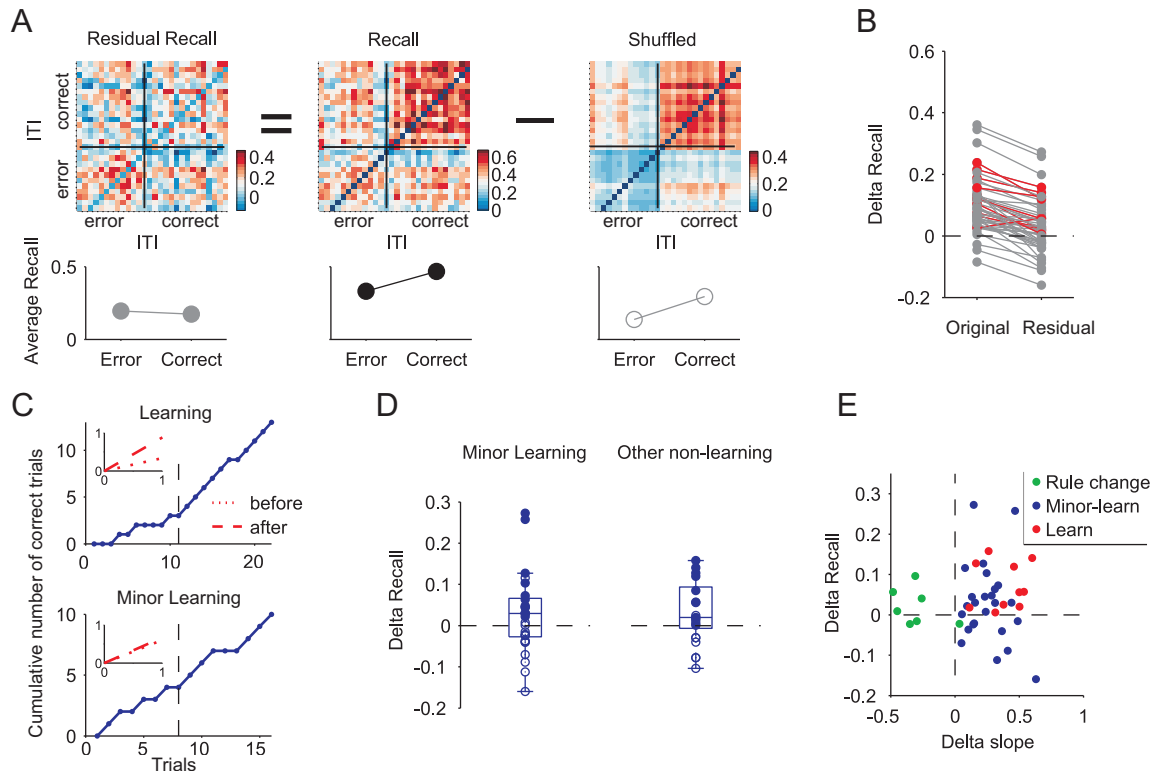
**Figure 2 - figure supplement 1.**

Statistics of neural populations and time periods during the inter-trial intervals for the 50 retained sessions.

(A) Proportion of neurons retained in the core ensemble were those that fired in every inter-trial interval. The black vertical dashed lines separate the sessions for each of the four rats.

(B) Durations of the inter-trial intervals within each session, given as the mean  $\pm$  SEM duration in seconds, separated into post-correct (filled symbols) and post-error (open symbols) inter-trial intervals. Red symbols are the learning sessions.

(C) Time spent along the maze during the inter-trial intervals, given as the mean  $\pm$  SEM seconds spent across all the animal and all the sessions for post-correct (black) and post-error (gray) inter-trial intervals. The maze has been linearised and divided in 5 equal sized sections, with position 1 being the reward location, position 3 the choice point of the Y-maze, and position 5 the end of the start arm – see Figure 5 (main text) for a schematic.



**Figure 2 - figure supplement 2.**

Recall of neural ensembles is learning-specific.

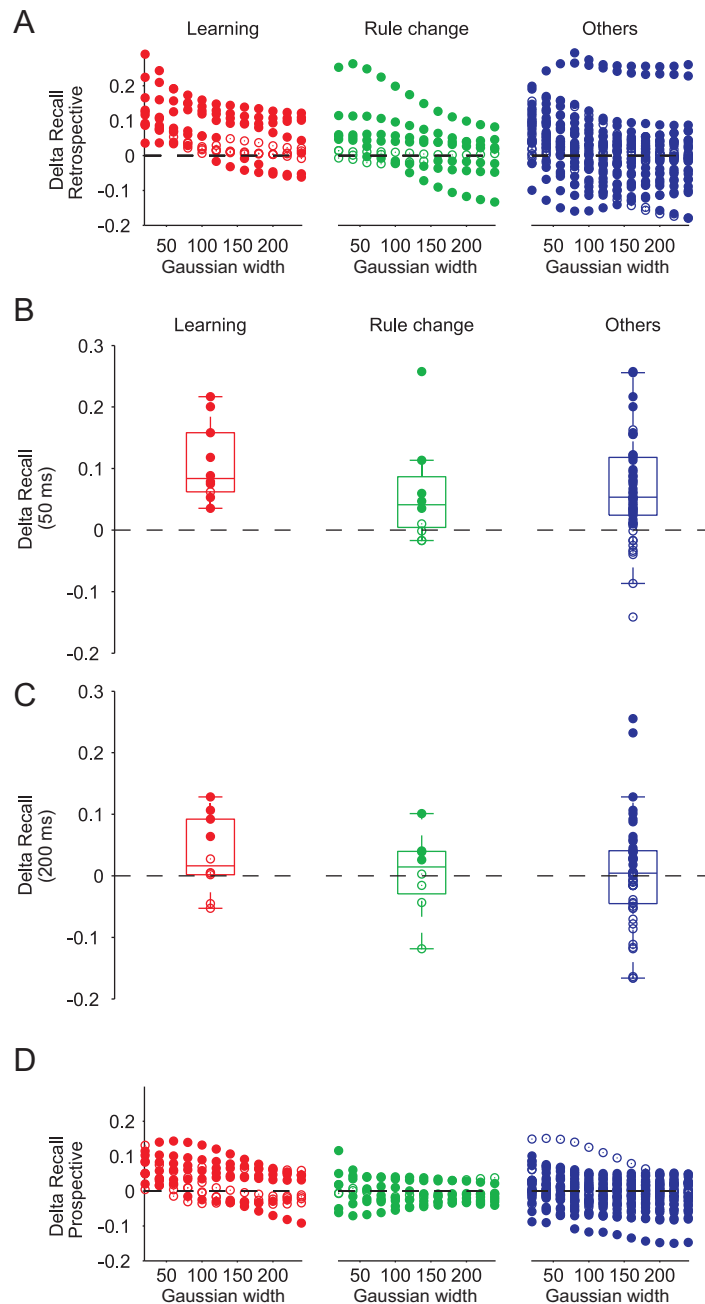
(A) Example correction of the recall matrix to remove the confounding effect of different durations of the post-correct and post-error inter-trial intervals. The residual recall matrix was obtained as the difference between the recall matrix and the mean matrix obtained from the shuffled inter-spike intervals (upper panels). For this example, the average recall values between error and correct intervals showed higher correlation among correct intervals in the shuffled model; this reversed the difference in recall between error and correct intervals (bottom panels) - in this case, ruling out a potentially higher recall during correct trials.

(B) Comparison of the difference between average correct recall and error recall (Delta recall) before and after correction by the shuffled control data. Red symbols are the learning sessions.

(C) To check whether the recall effect was specific to sessions showing abrupt learning (top panel; Figure 2, main text), we identified a subset of the other sessions with potential incremental or “minor” learning. These minor-learning sessions were any in which the curve of cumulative rewards contained a detectable upward inflection, as shown by the existence of any trial with a greater slope in a regression line after that trial than before it (insets, red lines). The vertical black dashed line is the identified learning trial.

(D) The difference between average correct recall and error recall (Delta recall) for the minor-learning and remaining other sessions. No systematic recall effect was observed for the minor-learning sessions, suggesting the recall effect was specific to abrupt transitions in behaviour.

(E) Relationship between behavioural change and the strength of recall. The difference between average correct recall and error recall (Delta recall) is plotted as a function of the difference between the slopes of the fitted lines before and after the learning trial (Delta slope). Sessions: learning (red), rule change (green), and minor-learning (blue). Delta slope for each rule change session was computed with respect to the rule change trial.



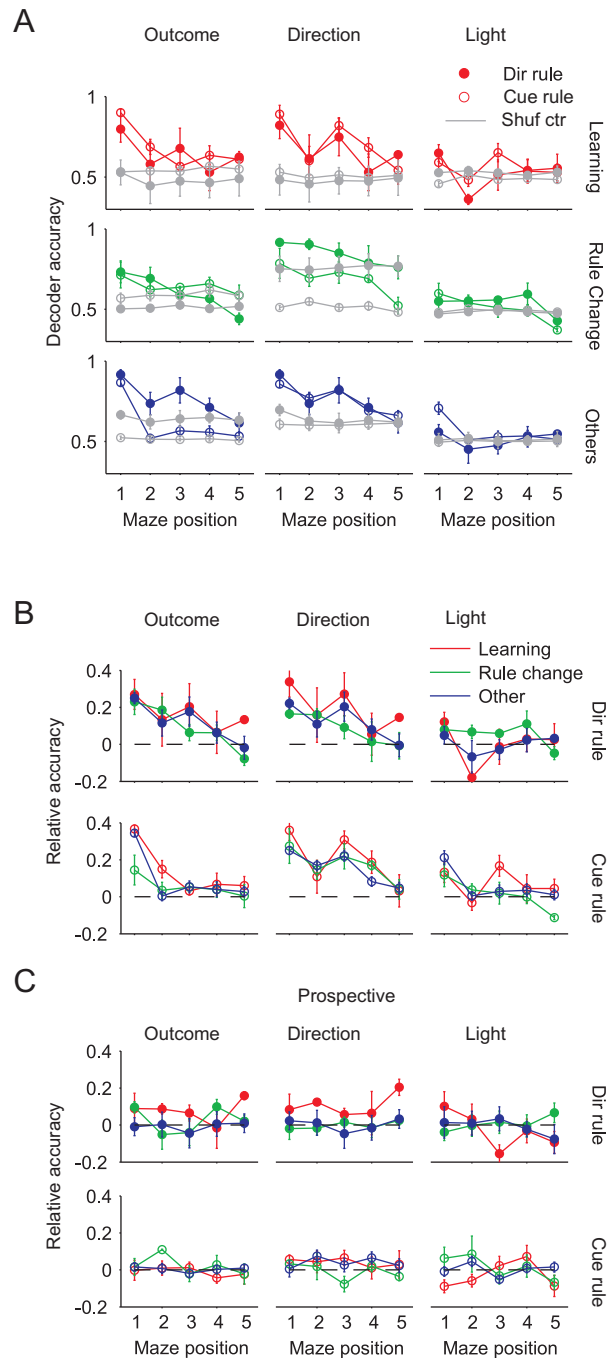
**Figure 2 - figure supplement 3.** Time-scale dependence of recall.

(A) Dependence of the recall of ensemble activity on the temporal precision of spike-train correlation. Here we plot the distribution of Delta recall across sessions as a function of the Gaussian width used to convolve the spike-trains. Retrospective recall, the difference in recall between intervals after correct and after error trials. Delta recall greater than zero indicates the interval similarity matrices were more correlated for correct than error intervals. Each symbol is one session. Filled circles indicate a difference at  $p < 0.05$  between the distributions of recall values in the error and correct intervals (Kolmogorov-Smirnov test).

(B) Comparison of Delta recall across learning, rule change, and other sessions after spike-train convolution with a Gaussian 50 ms wide.

(C) As for panel B, but for a Gaussian 200 ms wide.

(D) Dependence of the prospective recall of ensemble activity on the temporal precision of spike-train correlation. As for panel A, but here we plot the prospective Delta recall, the difference in recall between intervals before correct and before error trials. We only see a systematic recall at the smaller tested Gaussian widths ( $\leq 40$  ms), which is likely a reflection of the stronger retrospective recall effect at these widths.



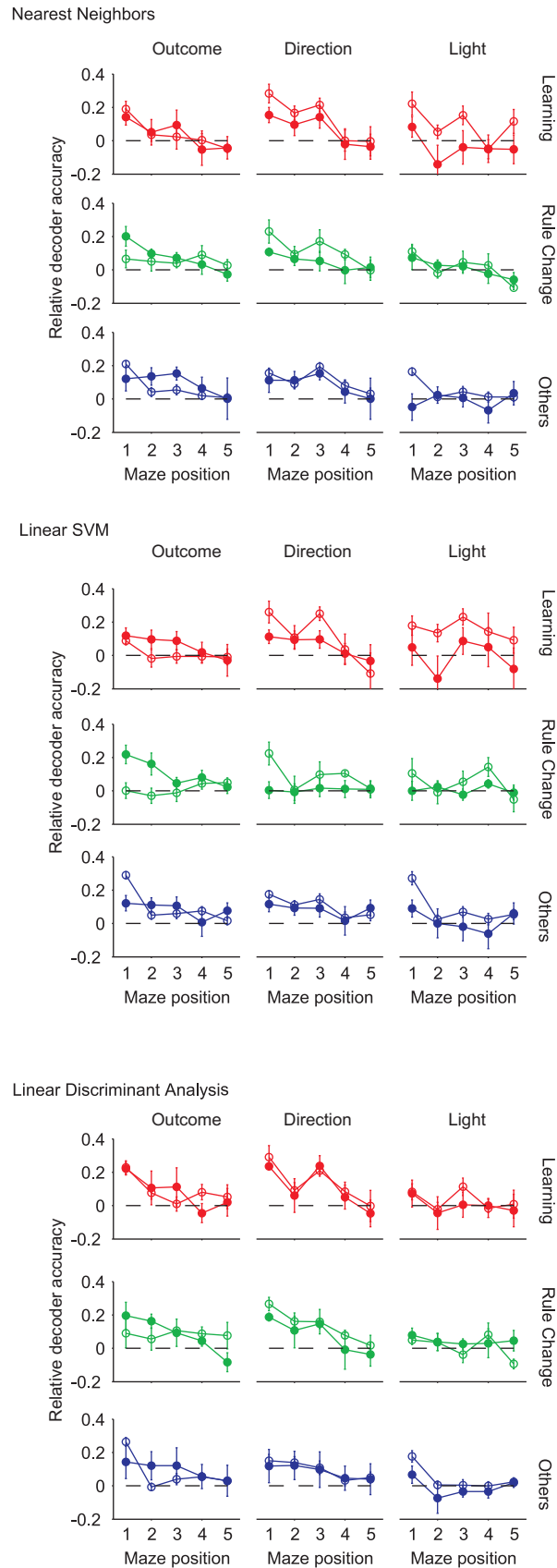
**Figure 5 - figure supplement 1.**

Further decoding analysis.

(A) Decoding can be near perfect. Here we replot the breakdown of the decoding results in Figure 5D as the absolute accuracy of the decoders (where 1 is maximum, indicating correct prediction of every held-out inter-trial interval). Each data point is the mean  $\pm$  SEM accuracy at that maze position. The control results on shuffled inter-trial interval labels are shown as grey lines. The “other” sessions were plotted in Figure 5C.

(B) Comparison of above-chance decoding performance between the same rule types. Each data point is the mean  $\pm$  SEM accuracy in excess of chance (dashed line) over the indicated combination of session type and rule type.

(C) Comparison of prospective decoding performance between the same rule types, confirming the absence of the prospective encoding of task-relevant information. Similar to panel B, here we plot the mean  $\pm$  SEM above-chance accuracy of decoding prospective outcome, direction, or cue position, separately for sessions with direction or light rules. As in panel B, decoding accuracy is normalised by the corresponding shuffled control decoding performance (where 0 is identical to shuffled controls).



**Figure 5 - figure supplement 2.**

Robustness of the retrospective encoding of task-relevant information. Using the same layout as Figure 5D, here we summarise the decoding performance of three further classifiers we tested on the data to check the robustness of the decoding results. Top: Nearest Neighbors; middle: linear Support Vector Machine; bottom: Linear Discriminant Analysis.