

---

*Research Articles: Behavioral/Cognitive*

## **Dissociable forms of uncertainty-driven representational change across the human brain**

**Matthew R. Nassar<sup>1</sup>, Joseph T. McGuire<sup>2</sup>, Harrison Ritz<sup>1</sup> and Joseph Kable<sup>3</sup>**

<sup>1</sup>*Department of Cognitive, Linguistic, and Psychological Sciences; Carney Institute for Brain Science, Brown University, Providence RI 02912-1821*

<sup>2</sup>*Department of Psychological & Brain Sciences; Boston University, Boston MA 02215*

<sup>3</sup>*Department of Psychology; University of Pennsylvania, Philadelphia PA 19143*

<https://doi.org/10.1523/JNEUROSCI.1713-18.2018>

Received: 7 July 2018

Revised: 7 November 2018

Accepted: 25 November 2018

Published: 6 December 2018

---

**Author contributions:** M.N., J.T.M., and J.K. designed research; M.N. and J.T.M. performed research; M.N. and H.R. analyzed data; M.N. wrote the first draft of the paper; M.N., J.T.M., H.R., and J.K. edited the paper; M.N. wrote the paper.

**Conflict of Interest:** The authors declare no competing financial interests.

We thank Ben Heasley for programming the task and Josh Gold, Michael J. Frank and Spencer Arbuckle for helpful discussion. This work was supported by NSF 1533623 and NIH R01-MH098899 to J.K., NSF BCS-1755757 and NIH F32-DA030870 to J.T.M., and NIH F32-MH102009-01A1 and NIH K99AG054732 to M.R.N.

Corresponding Author: Matthew R. Nassar, Department of Cognitive, Linguistic and Psychological Sciences, Brown University, Providence, RI 02912-1821, Phone: 607-316-4932, E-mail: [matthew\\_nassar@brown.edu](mailto:matthew_nassar@brown.edu)

**Cite as:** J. Neurosci 2018; 10.1523/JNEUROSCI.1713-18.2018

**Alerts:** Sign up at [www.jneurosci.org/alerts](http://www.jneurosci.org/alerts) to receive customized email alerts when the fully formatted version of this article is published.

Accepted manuscripts are peer-reviewed but have not been through the copyediting, formatting, or proofreading process.

Copyright © 2018 the authors

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46

**Dissociable forms of uncertainty-driven representational change across the human brain.**

Matthew R. Nassar<sup>1</sup>, Joseph T. McGuire<sup>2</sup>, Harrison Ritz<sup>1</sup> and Joseph Kable<sup>3</sup>

<sup>1</sup>Department of Cognitive, Linguistic, and Psychological Sciences; Carney Institute for Brain Science, Brown University, Providence RI 02912-1821

<sup>2</sup>Department of Psychological & Brain Sciences; Boston University, Boston MA 02215

<sup>3</sup>Department of Psychology; University of Pennsylvania, Philadelphia PA 19143

Abbreviated title: Dissociable explanations for representational change.

Figures: 3

Tables: 3

Abstract: 250

Introduction: 874

Discussion: 1228

Extended data: 1 figure

Corresponding Author:

Matthew R. Nassar

Department of Cognitive, Linguistic and Psychological Sciences

Brown University

Providence, RI 02912-1821

Phone: 607-316-4932

E-mail: [matthew\\_nassar@brown.edu](mailto:matthew_nassar@brown.edu)

**Acknowledgments:** We thank Ben Heasley for programming the task and Josh Gold, Michael J. Frank and Spencer Arbuckle for helpful discussion. This work was supported by NSF 1533623 and NIH R01-MH098899 to J.K., NSF BCS-1755757 and NIH F32-DA030870 to J.T.M., and NIH F32-MH102009-01A1 and NIH K99AG054732 to M.R.N.

**Abstract**

47 Environmental change can lead decision makers to shift rapidly among different  
48 behavioral regimes. These behavioral shifts can be accompanied by rapid changes in  
49 the firing pattern of neural networks. However, it is unknown what the populations  
50 of neurons that participate in such "network reset" phenomena are representing.  
51 Here we examined 1) whether and where rapid changes in multivariate activity  
52 patterns are observable with fMRI during periods of rapid behavioral change, and 2)  
53 what types of representations give rise to these phenomena. We did so by  
54 examining fluctuations in multi-voxel patterns of BOLD activity from male and  
55 female human subjects making sequential inferences about the state of a partially  
56 observable and discontinuously changing variable. We found that, within the  
57 context of this sequential inference task, the multivariate patterns of activity in a  
58 number of cortical regions contain representations that change more rapidly during  
59 periods of uncertainty following a change in behavioral context. In motor cortex,  
60 this phenomenon was indicative of discontinuous change in behavioral outputs,  
61 whereas in visual regions the same basic phenomenon was evoked by tracking of  
62 salient environmental changes. In most other cortical regions, including dorsolateral  
63 prefrontal and anterior cingulate cortex, the phenomenon was most consistent with  
64 directly encoding the degree of uncertainty. However, in a few other regions,  
65 including orbitofrontal cortex, the phenomenon was best explained by  
66 representations of a shifting context that evolve more rapidly during periods of  
67 rapid learning. These representations may provide a dynamic substrate for learning  
68 that facilitates rapid disengagement from learned responses during periods of  
69 change.

70

71

72

73

74

75 **Significance Statement**

76

77 Brain activity patterns tend to change more rapidly during periods of uncertainty

78 and behavioral adjustment, yet the computational role of such rapid transitions is  
79 poorly understood. Here we identify brain regions with fMRI BOLD activity patterns  
80 that change more rapidly during periods of behavioral adjustment and use  
81 computational modeling to attribute the phenomenon to specific causes. We  
82 demonstrate that the phenomenon emerges in different brain regions for different  
83 computational reasons, the most common being the representation of uncertainty  
84 itself, but that in a selective subset of regions including orbitofrontal cortex the  
85 phenomenon was best explained as a shifting latent state signal that may serve to  
86 control the degree to which recent temporal context affects ongoing expectations.

87

88

89

## 90 **Introduction**

91

92 Neural populations in rodent prefrontal cortex can undergo abrupt changes  
93 in firing concomitant with changes in performance in rule-based tasks (Durstewitz  
94 et al., 2010; Powell and Redish, 2016). Similar phenomena have been observed in  
95 the multi-voxel patterns in human fMRI data preceding changes in task strategy,  
96 leading to the notion that such changes might correspond to an “aha moment” at  
97 which the brain reorganizes to produce a new task set (Schuck et al., 2015). In  
98 rodent learning tasks that involve discontinuously changing reward contingencies,  
99 abrupt changes in firing of neurons in medial frontal cortex are observed more  
100 frequently during periods of uncertainty, during which animals appear to be  
101 searching for the best behavioral policy (Karlsson et al., 2012). It is unclear to what  
102 extent such phenomena are specific to medial frontal populations, or to what extent  
103 they might have an analog in human learning. Furthermore, while these “network  
104 resets” during periods of uncertainty are thought to play a role in behavioral  
105 flexibility in changing environments (Tervo et al., 2014) the exact computational  
106 role of abrupt changes in such neural representations remains unknown.

107 A number of different computational factors could explain previously  
108 observed network reset phenomena. First, and most simply, such abrupt changes

109 would be expected in a neural representation of the current behavioral policy,  
110 which in some cases may be directly related to the motor program. Successful  
111 execution of learning requires maintenance and updating of a behavioral policy,  
112 which would tend to change more rapidly during periods of uncertainty.

113         Alternatively, reset phenomena might result from representation of higher-  
114 order computational variables used to appropriately calibrate the rate of learning.  
115 Recent work has highlighted a number of computational variables that are  
116 important for successful learning in the presence of discontinuous environmental  
117 changes (change points). In particular, humans tend to increase rates of learning  
118 according to the probability with which a given outcome reflects a change point in  
119 the behavioral contingency (*change-point probability*) and according to the relative  
120 imprecision of their estimate of the current contingency (*relative uncertainty*)  
121 (Nassar et al., 2010; 2012). These computational variables both increase following  
122 change-points, albeit with different dynamics, to mediate rapid incorporation of  
123 new information during and after periods of environmental change. Change-point  
124 probability and relative uncertainty correlate with BOLD responses across a wide  
125 swath of brain regions including some that jointly reflect both variables and some  
126 that uniquely reflect either change-point probability or uncertainty (McGuire et al.,  
127 2014). In principle, neural representations of either computational factor might  
128 involve patterns of activation that mimic “network reset” phenomena, yet this  
129 possibility has never been tested directly.

130         Another signal that might give rise to reset-like dynamics is a continuously  
131 evolving latent state representation. Latent states, which represent the relevant  
132 behavioral context in cases where it is not directly observable, can improve learning  
133 in the face of abstract stimulus categories or repeated episodes by efficiently  
134 partitioning learning across distinct behaviorally relevant contexts (Gershman and  
135 Niv, 2010). While previous work has focused primarily on the advantage of such  
136 representations for rapid reinstatement of previously learned behaviors (Gershman  
137 et al., 2010; Wilson et al., 2014), another advantage of such representations is that  
138 they could facilitate rapid disengagement from established behaviors that are no  
139 longer relevant. By appropriately partitioning data collected over time in a changing

140 environment, such a mechanism could aid learning even if previously encountered  
141 environmental states do not recur. To accomplish this, such a latent state  
142 representation would need to evolve faster after a period of environmental change  
143 in order to effectively disengage from the previous behavioral context (Prescott  
144 Adams and MacKay, 2007; Wilson et al., 2010). While previous work has suggested  
145 that orbitofrontal cortex (OFC) might represent latent task states (Klein-Flugge et  
146 al., 2013; Stalnaker et al., 2014; Howard et al., 2015; Schuck et al., 2016; Howard and  
147 Kahnt, 2018) it is unclear whether such representations transition dynamically  
148 during periods of rapid learning as would be necessary to efficiently mediate  
149 disengagement of learned responses that are rendered irrelevant by environmental  
150 change.

151         Here we examined whether and where uncertainty-linked network resets are  
152 observable in human fMRI data, and evaluated the most likely computational  
153 explanation for these phenomena in individual brain regions. We did so using a  
154 multistep approach. First, we identified signals that change rapidly from trial to trial  
155 during periods of uncertainty and rapid learning and potentially correspond to  
156 network resets (Karlsson et al., 2012). Second, we generalized this notion of  
157 representational change across pairs of non-consecutive trials using  
158 representational similarity analysis (RSA) (Nili et al., 2014). Third, we formalized a  
159 set of candidate computational explanations for network-reset phenomena and  
160 allowed these explanations to compete to explain multivariate brain activity (Kragel  
161 et al., 2018).

162         We observed rapid changes in multivariate activity patterns across  
163 widespread cortical regions during periods of uncertainty and rapid learning. Using  
164 RSA, we showed that patterns in motor regions were best described as reflecting  
165 behavioral policy, patterns of activation in occipital regions were best described as  
166 registering the occurrence of change-points, and patterns across much of the rest of  
167 the cortex appeared to reflect uncertainty. However, patterns of activation in a small  
168 number of regions including OFC were most consistent with dynamic latent state  
169 representations, suggesting a possible role for the OFC in translating learning

170 signals into state changes that effectively disengage from behaviors learned in  
171 contexts that are no longer relevant.

172

173

## 174 **Methods**

175

### 176 *Behavioral task and analysis*

177 32 human subjects (17 female, 15 male) performed a computerized predictive  
178 inference task in an MRI scanner while undergoing functional neuroimaging. On  
179 each trial subjects were required to move a bucket to a new location at some point  
180 on the horizontal axis of a screen using a joystick controlled by the right hand and  
181 starting from a "home position" at the right-hand edge of the screen. Subjects were  
182 instructed to move the bucket to the inferred position of a helicopter, which was  
183 occluded by clouds, and thus not directly observable. Subjects had three seconds to  
184 place the bucket in their preferred location, after which the helicopter would drop a  
185 bag that contained either high value or neutral items (value designated by color,  
186 animation of falling bag lasted 1 second). The primary information in the task was  
187 provided by the horizontal location of the bag, which was selected at random on  
188 each trial from a normal distribution centered on the true helicopter location  
189 (incentivizing bucket placement under the inferred helicopter location) and with a  
190 standard deviation that was manipulated blockwise. Subjects completed four task  
191 blocks of 120 trials each (2 blocks of high/low standard deviation). On the majority  
192 of trials (90%) the helicopter would remain in the same location as in the previous  
193 trial, but occasionally (10%) the helicopter would relocate to a new position along  
194 the horizontal axis of the screen (selected randomly and uniformly). Since subjects  
195 could not see the helicopter, they were forced to infer its position based on history  
196 of previous bag locations, and in some cases were required to recognize and  
197 respond to a change in helicopter location. A more in-depth description of the  
198 behavioral task and an extensive characterization of subject behavior are available  
199 in our previous report (McGuire et al., 2014).

200

201 *MRI data acquisition and preprocessing*

202 T1-weighted MPAGE structural images (0.9375 X 0.9375 X 1mm voxels, 192  
203 X 256 matrix, 160 axial slices, TI=1100ms, TR=1630ms, TE=3.11ms, flip angle=15°),  
204 T2\*-weighted EPI functional data (3mm isotropic voxels, 64 X 64 matrix, 42 axial  
205 slices tilted 30° from the AC-PC plane, TR=2500ms, TE=25ms, flip angle=75°), and  
206 fieldmap images (TR=1000ms, TE=2.69 and 5.27ms, flip angle=60°) were acquired  
207 on a 3T Siemens Trio with a 32 channel head coil. Functional data were acquired in  
208 4 runs, each of which lasted 9 minutes and 25 seconds (226 images).

209 Data were preprocessed using AFNI (Cox, 1996; 2012) and FSL (Jenkinson et  
210 al., 2002; Smith et al., 2004; Jenkinson et al., 2012) in the following steps: 1) slice  
211 timing correction (AFNI's *3dTshift*), 2) motion correction (FSL's *MCFLIRT*), 3)  
212 fieldmap-based geometric undistortion, alignment with structural images, and  
213 registration to the MNI template (FSL's *FLIRT* and *FNIRT*), 4) spatial smoothing with  
214 a 6mm FWHM Gaussian kernel (FSL's *fslmaths*), 5) outlier attenuation (AFNI's  
215 *3dDespike*), and intensity-scaling by a single grand-mean value in each run (FSL's  
216 *fslmaths*). The resulting functional time series was deconvolved to estimate trial  
217 activations at the time of the bag drop using the least squares-separate method  
218 (Mumford et al., 2012) implemented in Matlab with inclusion of six rigid body  
219 motion parameters and sixteen low order cosine components (four per run) as  
220 regressors of no interest. Our decision to model the bag drop time point (as in our  
221 previous reports; (Nassar et al., 2012; McGuire et al., 2014)) was motivated by our  
222 interest in how the bag locations would affect internal representations. In practice,  
223 however, the rapid nature of our task prohibits us from making strong claims  
224 regarding the specificity of our results to a given task phase.

225 Alternative preprocessing pipelines were also used to verify the robustness  
226 of our findings (Tables 2&3). In one such pipeline the spatial smoothing was omitted  
227 from the pipeline described above, and instead spatial smoothing with a 6mm  
228 FWHM Gaussian kernel was applied to the coefficient maps resulting from  
229 representational similarity analysis. Another alternative preprocessing strategy  
230 omitted spatial smoothing and also implemented spatial pre-whitening to  
231 emphasize high frequency components of the spatial patterns (Walther et al., 2016).



232

233 *Computing normative dynamic learning rates*

234 Successful task performance required contending with imperfect cues about  
235 helicopter location (variability in the distribution of bag locations) as well as  
236 changes in helicopter location, which rendered past bag locations irrelevant to  
237 future ones. Optimal inference under such conditions can be achieved by applying  
238 Bayes rule to maintain and update a probability distribution over potential periods  
239 of stability, or run length (Prescott Adams and MacKay, 2007; Wilson et al., 2010).  
240 This solution can be approximated by using a single representative value for the run  
241 length, rather than maintaining the full distribution, yielding an error-driven  
242 learning rule in which the learning rate is adjusted dynamically from trial-to-trial:

243

$$B_{t+1} = B_t + \alpha_t \delta_t$$

244

245 where  $B_t$  reflects the underlying belief about helicopter location on trial  $t$ ,  $\delta$  reflects  
246 the prediction error on trial  $t$  (Difference between bag location and belief), and  $\alpha$   
247 reflects a dynamic learning rate, which varies from trial-to-trial and controls the  
248 influence of prediction errors on updated beliefs (Nassar et al., 2010; 2012).

249 Dynamic learning rates depend on change-point probability, or the  
250 probability that the helicopter relocated since the previous outcome was observed,  
251 and relative uncertainty, which reflects the fraction of uncertainty over the position  
252 of the upcoming bag location that is attributable to uncertainty about the current  
253 helicopter position (see Figure 1c; (Nassar et al., 2016)):

$$\alpha_t = \Omega_t + \tau_t - \Omega_t \tau_t$$

254

255 where  $\Omega_t$  represents change-point probability and  $\tau_t$  represents the relative  
256 uncertainty on trial  $t$ .

257 These latent variables were updated on each trial with a parameter-free  
258 normative model that took subject prediction errors as an input according to the  
259 following set of recursive equations:

260

$$\sigma_{\mu}^2 = \Omega_t \sigma_N^2 + (1 - \Omega_t) \sigma_N^2 \tau_t + \Omega_t (1 - \Omega_t) (\delta_t (1 - \tau_t))^2$$

261

$$\text{Relative uncertainty} = \tau_{t+1} = \frac{\sigma_{\mu}^2}{\sigma_{\mu}^2 + \sigma_N^2}$$

262

$$\text{Change point probability} = \Omega_{t+1} = \frac{\frac{H}{w}}{\frac{H}{w} + \mathcal{N}\left(\delta_{t+1} \mid 0, \frac{\sigma_N^2}{1 - \tau_{t+1}}\right)} (1 - H)$$

263

264

265 where  $\sigma_{\mu}^2$  is the total variance in beliefs about the helicopter location (the generative  
 266 mean),  $\sigma_N^2$  is the variance in the distribution of outcomes (bag drops) around that  
 267 mean,  $\delta_t$  is the prediction error, and  $H$  is the hazard rate and  $w$  is the width of the  
 268 screen. For a full derivation of the model and terms see (Nassar et al., 2010) and for  
 269 a complete description of the method for estimating latent variables see (Nassar et  
 270 al., 2016).

### 271 *Multivariate fMRI analysis*

272 Multivariate analyses were conducted in spherical searchlights (radius = 3  
 273 voxels) across the entire brain. Within each searchlight, the neural dissimilarity  
 274 between each pair of trials was computed as one minus the spatial Pearson  
 275 correlation between the voxel-wise activations for those trials.

276

### 277 *Trial-to-trial dissimilarity analysis*

278 Trial-to-trial dissimilarity scores were extracted by extracting the  $i=j-1$   
 279 diagonal elements from the dissimilarity matrix, which corresponded to the  
 280 dissimilarity between adjacent trials (see Figure 1d). The dissimilarity scores were  
 281 regressed onto an explanatory matrix containing an intercept, and dynamic learning  
 282 rates prescribed by a normative learning model, yielding one coefficient of interest  
 283 per subject, per searchlight.

284 As described above, dynamic learning rates depended on two factors:  
 285 change-point probability and relative uncertainty. In general, change-point

286 probability and relative uncertainty were both increased after change-points, albeit  
287 with different latencies, leading to learning rates that decay slowly as a function of  
288 time within context. Learning rates quantifying sensitivity to information provided  
289 on trial  $j$  were aligned with the trial-to-trial dissimilarity between trials  $j$  and  $j+1$ .  
290 Thus, our analysis targeted patterns of activity whose degree of change between  
291 trials  $j$  and  $j+1$  reflected normative learning predicted to occur from the outcome  
292 presented on trial  $j$ . The first 3 trials from each block were removed from analysis as  
293 they occurred at the onset of fMRI acquisition. Correlations between model-derived  
294 quantities (change-point probability, relative uncertainty, normative learning rate)  
295 and the six rigid body motion parameters (estimated from MCFLIRT and  
296 deconvolved using the least squares-separate method as described above) were  
297 uniformly small (mean Pearson  $R^2$  across participants  $< 0.009$  for each of the 18  
298 pairwise correlations) as were correlations with absolute relative displacements in  
299 the same measures (mean Pearson  $R^2$  across participants  $< 0.008$  for each of the 18  
300 pairwise correlations).

301

### 302 *Representational similarity analysis (RSA)*

303 Trial-to-trial dissimilarity analysis described above could be thought of as a  
304 special case of the general idea that the similarity between each pair of trials might  
305 be inversely related to the learning done between them. Because this pattern of  
306 similarity is what might be expected to emerge from a representation of the latent  
307 task state, which transitions abruptly from one context to the next and remains  
308 relatively stable after many trials in a well learned context, we will refer to it as the  
309 shifting latent state dissimilarity matrix. The hypothesis matrix for shifting latent  
310 states was generated by computing the extent to which the inference on trial  $i$  would  
311 factor into the inference on trial  $j$ , assuming normative learning:

312

$$H_{i,j} = 1 - \prod_{t=i}^{j-1} 1 - \alpha_t$$

313 where  $H$  is the shifting latent state dissimilarity matrix and  $\alpha$  is the learning rate  
314 prescribed by a normative model (Nassar et al., 2010), such that more prescribed  
315 learning between two trials corresponded to higher values of  $\alpha$ , a smaller product  
316 term, and thus a greater dissimilarity. The  $i=j-1$  diagonal of this matrix is  $1-(1-\alpha)_t$ , or  
317 just  $\alpha_t$ , and thus equivalent to the vector of trial-to-trial dissimilarities described  
318 above. However, the shifting latent state hypothesis matrix also includes  
319 information about other elements in the matrix, potentially offering a more  
320 powerful construct to ask a similar question. We examined whether this similarity  
321 structure was reflected in the neural dissimilarity between trials in each spherical  
322 searchlight. The lower triangle of the neural dissimilarity matrix was regressed onto  
323 a hypothesis matrix that included an intercept, the shifting latent state hypothesis  
324 matrix (lower triangle), and 15 dummy variables designed to remove the influence  
325 of autocorrelation on the coefficient of interest. These autocorrelation terms were  
326 derived from 15 off-diagonal binary matrices in which a single off diagonal ( $i = j-1$ ;  $i$   
327  $= j-2$ ;  $i = j-3$ ...  $i = j-15$ ) was set to one. These matrices were constructed to account  
328 for any variance in the neural dissimilarity matrices that could be explained by a  
329 fixed signal autocorrelation. To be sure that autocorrelation could not affect our  
330 analysis of interest, we also set all elements of the shifting latent state similarity  
331 matrix that fell outside of this range (trials separated by more than fifteen trials) to  
332 the maximum dissimilarity value.

333

#### 334 *Dissociating computational explanations with RSA*

335 To better understand the computations that give rise to rapid changes in  
336 neural patterns during periods of learning after a helicopter relocation, we  
337 constructed an exhaustive set of hypothesis matrices and conducted a  
338 representational similarity analysis in which these representations could compete  
339 with the shifting latent state matrix described above to explain structure in neural  
340 dissimilarity matrices. Thus, this analysis included the same shifting latent state  
341 matrix, but also included hypothesis matrices for various factors that could relate to  
342 task uncertainty, learning, or explain nuisance variance in the neural dissimilarity  
343 matrices. Hypothesis matrices were generated for three additional explanatory

344 variables of interest: 1) subject prediction (behavioral policy), 2) relative  
345 uncertainty, 3) change-point probability. We also included six additional nuisance  
346 variables: 4) the bag drop's location, 5) signed prediction error (ie, the distance  
347 between the prediction and the bag drop), 6) high CPP [to account for patterns of  
348 activity that may asymmetrically encode CPP], 7) high RU [to account for patterns of  
349 activity that may asymmetrically encode RU], 8) outcome reward value, and 9) task  
350 block. For factors 1-5 and 8, element (i,j) of the hypothesis matrix corresponded to  
351 the absolute difference in that factor on trials i and j. For factor 9, dissimilarity  
352 values were set to 0 for trials in the same block and 1 for trials in different blocks.  
353 Dissimilarity matrices for factors 6 & 7 were computed as one minus the  
354 multiplicative interaction of the model variable (6=change-point probability,  
355 7=relative uncertainty) on trials i and j, such that similarity was only hypothesized  
356 when the model-derived term took on a high value on both trials. These terms  
357 allowed the model to capture asymmetric representations of the two factors  
358 governing learning in our model, such as a representation that converged for values  
359 of high relative uncertainty but did not show any consistent pattern of activation  
360 when relative uncertainty was low.

361         The lower triangle of the neural dissimilarity matrix was extracted and  
362 regressed onto an explanatory matrix consisting of an intercept and the lower  
363 triangle of all hypothesis/nuisance matrices (including the shifting latent state and  
364 nuisance autocorrelation terms), yielding one coefficient per variable, per subject,  
365 per searchlight (Chikazoe et al., 2014; Kragel et al., 2018). Group level analyses were  
366 conducted by computing t-statistics across subjects for each variable and  
367 searchlight. Cluster-based permutation testing using cluster mass with a cluster  
368 forming threshold of  $p < 0.001$  and an alpha of 0.01 was used to identify significant  
369 activations (FSL's *randomise*) (Nichols and Holmes, 2002).

370

## 371 **Results**

372         To examine how neural signals change during periods of uncertainty we re-  
373 analyzed data from a previously published study that included recordings of fMRI  
374 BOLD signal and behavioral responses of human participants in a predictive

375 inference task (McGuire et al., 2014). Participants played a video game in which they  
376 tried to get as many coins as possible (redeemable for money) by catching bags of  
377 coins dropped from a hidden helicopter in the sky. Thus, on each task trial,  
378 participants estimated the state of an unobservable variable (the position of a  
379 helicopter) based on the history of an observable variable (the position of bags  
380 dropped from that helicopter) (McGuire et al., 2014). The task included abrupt  
381 change points at which the position of the helicopter was resampled from a uniform  
382 distribution, which forced participants to rapidly revise beliefs about the helicopter  
383 location in order to maintain successful task performance. Here we refer to periods  
384 of consistent helicopter position as contexts (Fig 1a), such that the task could be  
385 described as requiring dynamic belief updating both within (Fig 1a; vertical) and  
386 across (Fig 1a; horizontal) contexts.

387         As we described in our previous report, adjustments in the rate at which  
388 participants revised beliefs in response to new information were well described by a  
389 normative learning model that adjusted learning according to two computational  
390 variables: change-point probability and relative uncertainty (Fig 1b, compare pink  
391 and green lines; (McGuire et al., 2014; Nassar et al., 2016)). Change-point  
392 probability reflects the Bayesian posterior probability that the helicopter has  
393 relocated on the current trial, and is largest on trials with large spatial prediction  
394 errors (Fig 1c, blue line). Relative uncertainty captures the degree to which  
395 uncertainty about the true helicopter location should drive learning, is greatest on  
396 the trial after a spike in change-point probability, and decays as a function of trials  
397 thereafter (Fig 1c, yellow line). Both of these factors affect the sensitivity of ongoing  
398 beliefs to new information (e.g., bag locations), which can be expressed in terms of a  
399 dynamic learning rate (Fig 1c, green). We sought to identify relationships between  
400 the sensitivity of behavior to incoming information (i.e., learning rate) and the  
401 sensitivity of neural representations to the same information.

402

#### 403 *Representations change rapidly during learning*

404         The trial-to-trial dissimilarity in multivariate voxel activation patterns was  
405 related to the dynamic learning rates prescribed by the normative model (Fig 1d).

406 Trial wise neural dissimilarity was computed for each pair of sequentially adjacent  
407 trials using a whole brain searchlight procedure and regressed onto an explanatory  
408 matrix that included model-based estimates of dynamic learning rates. A  
409 constellation of regions showed patterns of activation that changed more rapidly  
410 during periods of rapid learning after change points (Fig 1e). These regions included  
411 OFC, but also clusters in dorsomedial frontal cortex (DMFC), occipital cortex, and the  
412 temporal lobe. Thus, with a simple measure of representational change, we  
413 identified neural signals whose representations updated more rapidly during  
414 periods of learning in multiple brain regions (cf. (Karlsson et al., 2012)).

415

416 *Testing for shifting latent state representations using RSA*

417 We next exploited representational similarity analysis (RSA) to extend and  
418 generalize the analysis above by incorporating information about the pairwise  
419 dissimilarity for all pairs of trials, not merely adjacent trial pairs. We hypothesized  
420 that the dissimilarity in neural representation for any pair of trials would depend on  
421 the cumulative amount of learning expected to occur between them under the  
422 normative model (see Methods). The hypothesized pattern of dissimilarity across  
423 trials is equivalent to what we would expect from a latent state representation that  
424 shifted rapidly at abrupt context transitions and concomitant periods of rapid  
425 learning, but remained relatively stable in periods when the statistics of the  
426 environment were stationary (Fig 2a). The pattern of dissimilarities predicted  
427 across adjacent trials using this strategy is exactly equivalent to the learning rates  
428 that served as the explanatory variable in the previous analysis (Fig 2b), but this  
429 generalization also makes predictions about the pattern of dissimilarities that would  
430 be observed across non-adjacent trials (Fig 2c). We used a searchlight to identify  
431 brain regions in which the neural dissimilarity matrix was positively associated with  
432 this hypothetical “shifting state representation” hypothesis matrix while controlling  
433 for fixed autocorrelation in the similarity structure (see Methods). A significant  
434 association was observed in a set of regions that overlapped with the results from  
435 the trial-wise dissimilarity analysis, including clusters in OFC, DMFC, occipital, and  
436 temporal regions (Fig 2d). As might be expected by the increased power owing to

437 the non-adjacent trial comparisons afforded by RSA analysis, we also identified  
438 additional regions that were not clearly indicated by our previous analysis including  
439 a number of visual regions, left motor cortex, and bilateral hippocampus (Fig 2d).

440

441 *Distinguishing between computational explanations for representational change*

442 We next sought to arbitrate among multiple possible causes for the varying  
443 rates of representational change. The rapid evolution of neural representations after  
444 change points might reflect different underlying computations in different brain  
445 regions. Our analysis focused on four candidate computations that could all  
446 theoretically drive network reset-like phenomena.

447 First, we considered the possibility that a brain region might reflect the  
448 behavioral policy of the participant. In our experimental task, the behavioral policy  
449 was reported directly by positioning a bucket at the predicted location (using a  
450 joystick) on each trial. For a given helicopter position, participants tended to place  
451 the bucket in a similar location, but changes in helicopter location corresponded to  
452 large changes in the bucket placement, which would correspond to abrupt  
453 transitions in a representation of behavioral policy after change points (Fig 3a).  
454 Occasionally, a new helicopter position was similar to one that had previously been  
455 encountered, such that a similar behavioral policy might be employed in two  
456 temporally separated contexts (Fig 3a; contexts 1&3).

457 A second possible explanation for rapid representational change after change  
458 points is that the representations could reflect the current level of change-point  
459 probability or relative uncertainty. Change-point probability changes most  
460 dramatically at a change in the context (Fig 1c), leading to predicted trialwise neural  
461 dissimilarity time courses that do the same (Fig 3b). The level of relative uncertainty  
462 changes most rapidly immediately after change-points (Figure 1c), and a neural  
463 representation of relative uncertainty should do the same (Fig 3c). However, either  
464 of these representations should return to a fixed pattern for all epochs across the  
465 experimental session that share the same level of change-point probability or  
466 relative uncertainty, irrespective of the current helicopter position (Fig 3b-c).



467           A final computational explanation for rapid representational changes after  
468 change points is that such a signal may reflect a latent state that is used to partition  
469 learning across distinct contexts (Wilson et al., 2014). For example, each new  
470 helicopter position could be reasonably thought of as a new temporal context,  
471 during which learning from prior contexts should be discounted to minimize  
472 interference (Fig 1a). Since the helicopter position cannot be resolved exactly, such  
473 a context representation would be expected to evolve over time in proportion to the  
474 rate of learning about the current context. This idea was formalized in figure 2, and  
475 as described previously, would lead to latent state representations that change  
476 rapidly at change points and immediately afterwards and change only minimally  
477 during periods of prolonged stability (Fig 3d). Unlike the other computational  
478 factors discussed above, a latent state representation would not necessarily exhibit  
479 any systematic similarity relation between one context and another – as our task did  
480 not include situations in which the helicopter returned exactly to a previously  
481 occupied position. Such a latent state signal might provide an evolving substrate to  
482 which outcomes could be linked in order to achieve rational adjustments of  
483 learning.

484           Each of these representations would yield more rapid changes in neural  
485 patterns after change points in our task, and indeed, they make very similar  
486 predictions for how neural dissimilarity metrics between adjacent trials should  
487 evolve over time (Fig 3 middle column, top plots). Predictions of trial-to-trial  
488 dissimilarity made for the four candidate computations were highly correlated (all  
489 average pairwise Pearson correlations [ $r$ ] were greater than 0.45, with predictions  
490 for shifting latent representations particularly highly correlated with those for  
491 relative uncertainty [ $r = 0.80$ ] and behavioral policy [ $r = 0.74$ ]), suggesting that the  
492 representations of these computations could not be distinguished based on  
493 adjacent-trial dissimilarity alone.

494           However, the four candidate representations differed drastically in their  
495 predictions about the dissimilarity for non-adjacent pairs of trials. We constructed  
496 hypothesis matrices for each candidate representation by considering the expected  
497 difference in the computation of interest across all possible pairs of trials. These

498 hypothesis matrices highlight qualitative features of each candidate computation;  
499 behavioral policy frequently undergoes abrupt shifts but often takes on a similar  
500 value to a previous state, change-point probability highlights differences between  
501 change point and non-change point trials, relative uncertainty highlights the  
502 differences between high relative uncertainty and other trials, and shifting latent  
503 states capture differences largely near the diagonal (Fig 3, middle column, bottom).  
504 Consistent with these qualitative differences, correlations between the hypothesis  
505 matrices for the different candidate representations were relatively low (all  
506 pairwise  $r < 0.16$ ), suggesting that the candidate representations could be efficiently  
507 distinguished when considering the entire pairwise dissimilarity matrix.

508         We exploited these distinct predictions using a representational similarity  
509 analysis approach that allowed alternative explanations of representational change  
510 to compete to explain the observed neural dissimilarity matrix. Neural dissimilarity  
511 was computed for each pair of trials as one minus the spatial correlation of trial-  
512 activations across voxels in a searchlight and regressed onto an explanatory matrix  
513 that included the hypothesis matrices for all four candidate representations, along  
514 with a number of other explanatory terms designed to account for factors changing  
515 throughout the task and simple sources of variability such as autocorrelation (see  
516 Methods).

517         Representational similarity analysis supported distinct explanations for  
518 representational change in different anatomical regions. Behavioral policy provided  
519 a good description of BOLD activity patterns in left motor cortex (contralateral to  
520 the hand used to move the joystick and execute the behavioral policy) and visual  
521 cortex (Figure 3a, right; Table 1). Representations of change-point probability were  
522 prominent in occipital cortex and precuneus (Figure 3b; Table 1). Representations  
523 of relative uncertainty were widespread across the brain and included DMFC,  
524 dorsolateral prefrontal cortex, bilateral parietal cortices, insula, as well as some  
525 occipital and temporal regions (Figure 3c, right). Patterns of activation consistent  
526 with a latent state that shifts according to assessment of the current context were  
527 prominent in OFC and temporal cortex (Fig 3d, right; Table 1).

528           The relationship between the neural dissimilarity in OFC and the  
529 dissimilarity structure predicted by a shifting latent state signal was unlikely to be  
530 an artifact of motion or eye movements. Normative learning rate, the primary driver  
531 of the shifting latent state hypothesis matrix, was not correlated with motion  
532 parameters to any significant degree (average Pearson R2 across subjects < 0.006  
533 for each of the 6 motion parameters), nor was it correlated with eye-movements in a  
534 follow up study using the same task run outside of the scanner (McGuire et al.,  
535 2014).

536           Our findings were robust to analysis choices including those affecting the  
537 spatial frequency of our multivariate fMRI signals. There is active debate over best  
538 practices in pre-processing fMRI data for representational similarity analysis with  
539 some work supporting liberal spatial smoothing of raw data prior to analysis (de  
540 Beeck, 2010; Hendriks et al., 2017) and other work suggesting that excessive spatial  
541 smoothing can dampen signals of interest by reducing high frequency components  
542 of the signal (Gardumi et al., 2016). As we had no a priori predictions about the  
543 spatial scale of our signal, we repeated our full representational similarity analysis  
544 on unsmoothed fMRI data, instead adding an additional smoothing step after  
545 representational similarity analysis on the resulting coefficient maps. This analysis  
546 yielded very similar results to our original analysis (compare Fig 3-1 extended data  
547 to Fig 3), including similar shifting latent state effects in OFC (Fig 3-1d extended  
548 data, right; Table 2). An ROI analysis applied to peak voxels for the shifting latent  
549 state regressor in our primary analysis (Table 1) that further emphasized high  
550 frequency components of spatial pattern by using a pre-whitening procedure  
551 (Walther et al., 2016) confirmed that OFC latent state representations were evident  
552 even when neural dissimilarity was restricted to high spatial frequency information  
553 ( $p < 0.05$ ; Table 3).

554           The observed shifting latent state effects in OFC were not driven by  
555 relationships between additional explanatory variables included in the regression  
556 model, as exclusion of other explanatory variables yielded very similar relationships  
557 (Table 3). It is noteworthy that this was not true of all clusters that survived whole-  
558 brain correction in our representational similarity regression analysis; clusters

559 identified in left superior parietal lobule and right occipital cortex were not related  
560 to the shifting latent state predictions in isolation (Table 3). Furthermore, the  
561 relationship between shifting latent state predictions and OFC patterns of activation  
562 was also robust to our assumptions about the exact timing of learning; a time shifted  
563 version of the shifting latent state hypothesis matrix that assumed learning occurred  
564 immediately upon observing a trial outcome could also describe similarity patterns  
565 observed in right and left OFC (Table 3).

566 In summary, while we found a number of regions that showed rapidly  
567 changing representations during periods of uncertainty following a context change,  
568 these reset-like phenomena were due to dissociable computational explanations.  
569 While a few regions were implicated in representing behavioral policy or change-  
570 point probability, most of these regions reflected relative uncertainty, and a smaller  
571 subset of regions including OFC were consistent with representing a latent state that  
572 is adjusted according to changes in context.

573

## 574 **Discussion**

575

576 Neural representations in rodent medial frontal cortex rapidly change during  
577 periods of uncertainty (Karlsson et al., 2012). Here we demonstrate, in the context  
578 of a dynamic learning task, that such rapid representational changes are present in  
579 the BOLD signal in widespread cortical and subcortical regions. Furthermore, we  
580 showed that these rapid representational changes are consistent with several  
581 different computational explanations, which could be teased apart by considering  
582 the similarity structure of non-adjacent trials through representational similarity  
583 analysis.

584 Our analyses revealed distinct explanations for rapid representational  
585 changes in different brain regions. Focal representations of behavioral policy and  
586 change-point probability were identified in motor and visual cortex respectively,  
587 while widespread representations of relative uncertainty were observed throughout  
588 the brain. In addition, a small number of brain areas including the OFC had patterns

589 of activation consistent with a form of shifting latent state representation that could  
590 speed disengagement from well-learned responses in a changing context.

591       Perhaps most straightforwardly, our analysis revealed that left motor cortex  
592 contained representations consistent with behavioral policy. In our task, this policy  
593 was completely concordant with the physical movement necessary to implement  
594 the behavioral policy. Thus, we interpret these results as a consequence of our  
595 experimental design, which required subjects to provide an analog behavioral  
596 output of their behavioral policy with their right hand on each task trial. Thus, this  
597 result was likely driven, at least in part, by a univariate effect of movement  
598 magnitude in the contralateral motor cortex.

599       Two other computations that we identified using this approach, change-point  
600 probability and relative uncertainty, had been the focus of a previous paper using  
601 this same dataset (McGuire et al., 2014). In the case of change-point probability,  
602 both univariate and RSA analyses revealed occipital cortex and precuneus as the  
603 locus of neural representation (see Figure 2c and (McGuire et al., 2014)). However,  
604 relative uncertainty representations identified using RSA were considerably more  
605 widespread than those identified through univariate activations (see Figure 2c and  
606 (McGuire et al., 2014)). This broader set of areas included some regions that were  
607 activated in the univariate analysis (e.g., DMFC), some that were deactivated in the  
608 univariate analysis (e.g., ventromedial prefrontal cortex), and some that were not  
609 identified in univariate analyses at all (e.g., temporal cortex). The near-ubiquitous  
610 cortical representation of relative uncertainty revealed by RSA is somewhat  
611 analogous to the widespread representations of reward prediction errors that have  
612 been identified using multivariate fMRI analysis methods (Vickery et al., 2011).  
613 Interestingly, both reward prediction errors and relative uncertainty have been  
614 suggested to be signaled through brainstem neuromodulatory systems that could  
615 potentially have widespread effects throughout the brain (Schultz, 1997; Yu and  
616 Dayan, 2005; Doya, 2008; Nassar et al., 2012).

617       In addition to providing a more sensitive tool to identify well-specified  
618 computational variables, RSA also allowed us to look for patterns of activity that  
619 could not easily be detected in univariate analyses. In particular, it allowed us to

620 look for neural representations of a dynamically shifting state representation,  
621 without making strong assumptions about what the signal would look like at any  
622 given moment. It has been proposed that state representations provided by the OFC  
623 might serve to hasten learning in environments that include a small number of  
624 repeated contexts (Gershman and Niv, 2010; Wilson et al., 2014; Schuck et al.,  
625 2016). This proposal is supported by observations that OFC representations encode  
626 the predicted identities of action outcomes (Klein-Flugge et al., 2013; Stalnaker et  
627 al., 2014; Howard et al., 2015; Howard and Kahnt, 2018), can reflect a probability  
628 distribution over the causal source of outcomes (Chan et al., 2016), and can be used  
629 to decode latent states that control action-outcome contingency (Schuck et al.,  
630 2016). Here we hypothesized that shifts in the same type of latent state  
631 representations might implement the rapid learning that should and does follow  
632 change-points in outcome contingencies (Prescott Adams and MacKay, 2007; Nassar  
633 et al., 2010; Wilson et al., 2010). Such an implementation could make use of existing  
634 computational elements to efficiently partition learned associations that pertain to  
635 distinct and unrelated contexts, effectively creating the product partitions necessary  
636 for optimal inference amid change-points (Prescott Adams and MacKay, 2007).

637         In line with this idea, we identified signals in OFC consistent with a shifting  
638 state signal that changed more rapidly during periods of learning. The region of OFC  
639 that we identified included both lateral regions (Fig 3d), similar to those where  
640 outcome identity representations have previously been observed (Howard and  
641 Kahnt, 2018) and more medial regions (Fig 3d) closer to where state  
642 representations have previously been reported (Schuck et al., 2016). Nonetheless,  
643 the OFC regions in which we identify shifting state signals are still somewhat lateral  
644 to those reported by Schuck and colleagues, and future work should examine  
645 whether the sorts of abrupt transitions in representation that we identify here  
646 indeed occur the same regions that as those that seem to represent state within a  
647 cognitive map of task space.

648         A neural population that encoded such a shifting state signal would be well  
649 positioned to transform a direct representation of dynamic learning rate, such as  
650 have been identified in cortical regions (Behrens et al., 2007; Krugel et al., 2009;

651 McGuire et al., 2014) and thought to be broadcast through noradrenergic  
652 neuromodulation (Yu and Dayan, 2005; Nassar et al., 2012; Browning et al., 2015),  
653 into a proportional change in associative strength. Using a learning signal to control  
654 the rate of contextual shift could enable a simple associative neural network to  
655 accomplish the type of adaptive learning that has previously been modeled as a  
656 delta-rule update with a varying learning rate. In such a case, increases in apparent  
657 learning would be implemented through changes in the substrate for learning, or  
658 the active latent state, rather than by adjusting associative strength per se.

659         Representations of latent state that transition dynamically from one context  
660 to the next are similar in spirit to the concept of event segmentation in episodic  
661 memory (Ezzyat and Davachi, 2010). Segmenting events is useful in that it can allow  
662 memories that are embedded within the same event but separated in time to share  
663 associations, while memories that may be closer in time but embedded in separate  
664 events are maintained separately, preventing interference (Reynolds et al., 2007).  
665 One mechanism through which segmentation could be achieved involves dynamic  
666 adjustment of the time-constant in slowly fluctuating temporal context signals to  
667 effectively “reset” context at event boundaries (Howard and Kahana, 2002; Howard  
668 et al., 2010; Manning et al., 2011). Our data suggest a link between this aspect of  
669 episodic encoding and the dynamic adjustments of learning that have been observed  
670 at context boundaries (Behrens et al., 2007; Nassar et al., 2010; McGuire et al.,  
671 2014). However, aspects of our findings also raise questions about the extent of this  
672 link. While our results could be interpreted as supporting roles for OFC and  
673 temporal lobe in segmenting contexts, we did not observe the same phenomenon in  
674 the hippocampus, which is thought to play a key role in event segmentation (Ezzyat  
675 and Davachi, 2014; Hsieh et al., 2014; Shapiro, 2014). Instead, we found that  
676 representations in hippocampus, like many other brain regions, were best explained  
677 as representing uncertainty itself. One potentially relevant detail is that previous  
678 contexts were not systematically re-visited in our task, reducing demands for  
679 episodic retrieval. An interesting avenue for future work would be to examine how  
680 the representations we identified respond when the context abruptly returns to a

681 previously encountered state, such as might require a form of mental time travel for  
682 successful performance (Manning et al., 2011).

683 Our results, especially regarding the OFC, demonstrate the utility of  
684 analyzing the representational similarity of multi-voxel patterns of activity in  
685 concert with computational modeling. Such an approach allowed us to identify  
686 neural representations consistent with a specific computational role for OFC, which  
687 in principle could not have been isolated in our task with univariate activation or  
688 multivariate classification analyses.

689 In summary, we show that shifts in the statistics of the environment during a  
690 dynamic learning task induced both elevated learning and changes in neural  
691 representation. These changes in neural representation were attributed to specific  
692 computations using RSA. Our results identified widespread representations of  
693 relative uncertainty throughout the brain, together with more focal representations  
694 of change-point probability and behavioral policy. In addition, a small number of  
695 brain areas including the OFC had patterns of activation consistent with a shifting  
696 latent state representation that could speed unlearning of irrelevant information in  
697 a changing context.

698

699 **Reference:**

700

701 Behrens TEJ, Woolrich MW, Walton ME, Rushworth MFS (2007) Learning the value  
702 of information in an uncertain world. *Nature Neuroscience* 10:1214–1221.

703 Browning M, Behrens TE, Jocham G, O'Reilly JX, Bishop SJ (2015) Anxious  
704 individuals have difficulty learning the causal statistics of aversive  
705 environments. *Nature Neuroscience* 18:590–596.

706 Chan SCY, Niv Y, Norman KA (2016) A Probability Distribution over Latent Causes,  
707 in the Orbitofrontal Cortex. *Journal of Neuroscience* 36:7817–7828.

708 Chikazoe J, Lee DH, Kriegeskorte N, Anderson AK (2014) Population coding of affect  
709 across stimuli, modalities and individuals. *Nature Neuroscience* 17:1114–1122.

710 Cox RW (1996) AFNI: software for analysis and visualization of functional magnetic  
711 resonance neuroimages. *Comput Biomed Res* 29:162–173.

712 Cox RW (2012) AFNI: what a long strange trip it's been. *NeuroImage* 62:743–747.



- 713 de Beeck HPO (2010) Against hyperacuity in brain reading: Spatial smoothing does  
714 not hurt multivariate fMRI analyses? *NeuroImage* 49:1943–1948.
- 715 Doya K (2008) Modulators of decision making. *Nature Neuroscience* 11:410–416.
- 716 Durstewitz D, Vittoz NM, Floresco SB, Seamans JK (2010) Abrupt Transitions  
717 between Prefrontal Neural Ensemble States Accompany Behavioral Transitions  
718 during Rule Learning. *Neuron* 66:438–448.
- 719 Ezzyat Y, Davachi L (2010) What Constitutes an Episode in Episodic Memory?  
720 *Psychol Sci* 22:243–252.
- 721 Ezzyat Y, Davachi L (2014) Similarity Breeds Proximity: Pattern Similarity within  
722 and across Contexts Is Related to Later Mnemonic Judgments of Temporal  
723 Proximity. *Neuron* 81:1179–1189.
- 724 Gardumi A, Ivanov D, Hausfeld L, Valente G, Formisano E, Uludağ K (2016) The effect  
725 of spatial resolution on decoding accuracy in fMRI multivariate pattern analysis.  
726 *NeuroImage* 132:32–42.
- 727 Gershman SJ, Blei DM, Niv Y (2010) Context, learning, and extinction. *Psychological*  
728 *Review* 117:197–209.
- 729 Gershman SJ, Niv Y (2010) Learning latent structure: carving nature at its joints.  
730 *Current Opinion in Neurobiology* 20:251–256.
- 731 Hendriks MHA, Daniels N, Pegado F, Op de Beeck HP (2017) The Effect of Spatial  
732 Smoothing on Representational Similarity in a Simple Motor Paradigm. *Front*  
733 *Neurol* 8:222.
- 734 Howard JD, Gottfried JA, Tobler PN, Kahnt T (2015) Identity-specific coding of future  
735 rewards in the human orbitofrontal cortex. *Proceedings of the National*  
736 *Academy of Sciences* 112:5195–5200.
- 737 Howard JD, Kahnt T (2018) Identity prediction errors in the human midbrain  
738 update reward-identity expectations in the orbitofrontal cortex. *Nature*  
739 *Communications*:1–11.
- 740 Howard MW, Kahana MJ (2002) A Distributed Representation of Temporal Context.  
741 *Journal of Mathematical Psychology* 46:269–299.
- 742 Howard MW, Shankar KH, Jagadisan UKK (2010) Constructing Semantic  
743 Representations From a Gradually Changing Representation of Temporal  
744 Context. *Top Cogn Sci* 3:48–73.
- 745 Hsieh L-T, Gruber MJ, Jenkins LJ, Ranganath C (2014) Hippocampal Activity Patterns  
746 Carry Information about Objects in Temporal Context. *Neuron* 81:1165–1178.

- 747 Jenkinson M, Bannister P, Brady M, Smith S (2002) Improved optimization for the  
748 robust and accurate linear registration and motion correction of brain images.  
749 *NeuroImage* 17:825–841.
- 750 Jenkinson M, Beckmann CF, Behrens TEJ, Woolrich MW, Smith SM (2012) FSL.  
751 *NeuroImage* 62:782–790.
- 752 Karlsson MP, Tervo DGR, Karpova AY (2012) Network resets in medial prefrontal  
753 cortex mark the onset of behavioral uncertainty. *Science* 338:135–139.
- 754 Klein-Flugge MC, Barron HC, Brodersen KH, Dolan RJ, Behrens TEJ (2013)  
755 Segregated Encoding of Reward-Identity and Stimulus-Reward Associations in  
756 Human Orbitofrontal Cortex. *Journal of Neuroscience* 33:3202–3211.
- 757 Kragel PA, Kano M, Van Oudenhove L, Ly HG, Dupont P, Rubio A, Delon-Martin C,  
758 Bonaz BL, Manuck SB, Gianaros PJ, Ceko M, Reynolds Losin EA, Woo C-W,  
759 Nichols TE, Wager TD (2018) Generalizable representations of pain, cognitive  
760 control, and negative emotion in medial frontal cortex. *Nature Publishing Group*.
- 761 Krugel LK, Biele G, Mohr PNC, Li S-C, Heekeren HR (2009) Genetic variation in  
762 dopaminergic neuromodulation influences the ability to rapidly and flexibly  
763 adapt decisions. *Proceedings of the National Academy of Sciences* 106:17951–  
764 17956.
- 765 Manning JR, Polyn SM, Baltuch GH, Litt B, Kahana MJ (2011) Oscillatory patterns in  
766 temporal lobe reveal context reinstatement during memory search. *Proceedings*  
767 *of the National Academy of Sciences* 108:12893–12897.
- 768 McGuire JT, Nassar MR, Gold JI, Kable JW (2014) Functionally dissociable influences  
769 on learning rate in a dynamic environment. *Neuron* 84:870–881.
- 770 Mumford JA, Turner BO, Ashby FG, Poldrack RA (2012) Deconvolving BOLD  
771 activation in event-related designs for multivoxel pattern classification analyses.  
772 *NeuroImage* 59:2636–2643.
- 773 Nassar MR, Bruckner R, Gold JI, Li S-C, Heekeren HR, Eppinger B (2016) Age  
774 differences in learning emerge from an insufficient representation of  
775 uncertainty in older adults. *Nature Communications* 7:11609.
- 776 Nassar MR, Rumsey KM, Wilson RC, Parikh K, Heasley B, Gold JI (2012) Rational  
777 regulation of learning dynamics by pupil-linked arousal systems. *Nature*  
778 *Neuroscience* 15:1040–1046.
- 779 Nassar MR, Wilson RC, Heasley B, Gold JI (2010) An approximately Bayesian delta-  
780 rule model explains the dynamics of belief updating in a changing environment.  
781 *Journal of Neuroscience* 30:12366–12378.

- 782 Nichols TE, Holmes AP (2002) Nonparametric permutation tests for functional  
783 neuroimaging: a primer with examples. *Hum Brain Mapp* 15:1–25.
- 784 Nili H, Wingfield C, Walther A, Su L, Marslen-Wilson W, Kriegeskorte N (2014) A  
785 Toolbox for Representational Similarity Analysis. *PLoS Comput Biol*  
786 10:e1003553.
- 787 Powell NJ, Redish AD (2016) Representational changes of latent strategies in rat  
788 medial prefrontal cortex precede changes in behaviour. *Nature Communications*  
789 7:12830.
- 790 Prescott Adams R, MacKay DJC (2007) Bayesian Online Changepoint Detection.  
791 eprint arXiv:07103742:–.
- 792 Reynolds JR, Zacks JM, Braver TS (2007) A computational model of event  
793 segmentation from perceptual prediction. *Cogn Sci* 31:613–643.
- 794 Schuck NW, Cai MB, Wilson RC, Niv Y (2016) Human Orbitofrontal Cortex  
795 Represents a Cognitive Map of State Space. *Neuron* 91:1402–1412.
- 796 Schuck NW, Gaschler R, Wenke D, Heinzle J, Frensch PA, Haynes J-D, Reuber C  
797 (2015) Medial Prefrontal Cortex Predicts Internally Driven Strategy Shifts.  
798 *Neuron* 86:331–340.
- 799 Schultz W (1997) A Neural Substrate of Prediction and Reward. *Science* 275:1593–  
800 1599.
- 801 Shapiro ML (2014) Time and Again. *Neuron* 81:964–966.
- 802 Smith SM, Jenkinson M, Woolrich MW, Beckmann CF, Behrens TEJ, Johansen-Berg H,  
803 Bannister PR, De Luca M, Drobnjak I, Flitney DE, Niazy RK, Saunders J, Vickers J,  
804 Zhang Y, De Stefano N, Brady JM, Matthews PM (2004) Advances in functional  
805 and structural MR image analysis and implementation as FSL. *NeuroImage* 23  
806 Suppl 1:S208–S219.
- 807 Stalnaker TA, Cooch NK, McDannald MA, Liu T-L, Wied H, Schoenbaum G (2014)  
808 Orbitofrontal neurons infer the value and identity of predicted outcomes. *Nature*  
809 *Communications* 5:3926.
- 810 Tervo DGR, Proskurin M, Manakov M, Kabra M, Vollmer A, Branson K, Karpova AY  
811 (2014) Behavioral Variability through Stochastic Choice and Its Gating by  
812 Anterior Cingulate Cortex. *Cell* 159:21–32.
- 813 Vickery TJ, Chun MM, Lee D (2011) Ubiquity and Specificity of Reinforcement  
814 Signals throughout the Human Brain. *Neuron* 72:166–177.
- 815 Walther A, Nili H, Ejaz N, Alink A, Kriegeskorte N, Diedrichsen J (2016) Reliability of

816       dissimilarity measures for multi-voxel pattern analysis. *NeuroImage* 137:188–  
817       200.

818       Wilson RC, Nassar MR, Gold JI (2010) Bayesian online learning of the hazard rate in  
819       change-point problems. *Neural Comput* 22:2452–2476.

820       Wilson RC, Takahashi YK, Schoenbaum G, Niv Y (2014) Orbitofrontal cortex as a  
821       cognitive map of task space. *Neuron* 81:267–279.

822       Yu AJ, Dayan P (2005) Uncertainty, neuromodulation, and attention. *Neuron*  
823       46:681–692.

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846 Figure 1: **Trialwise neural dissimilarity is increased after change-points during periods of**  
847 **rapid learning for multiple brain regions. A)** Participants were asked to move a bucket (pink  
848 rectangle) on each trial to the location most likely to deliver a reward (in the form of a bag containing  
849 coins). On each trial (stacked vertically) the participant would observe an outcome (bag location;  
850 gray circle) that they could use to update their bucket placement for the subsequent trial (orange  
851 arrow). Most contiguous trials were generated from the same context, which was defined by a fixed  
852 outcome distribution, however at occasional change points, the context (mean outcome location)  
853 shifted abruptly and unpredictably. **B)** An example sequence of outcomes (gray circles) and  
854 corresponding participant bucket placements (pink line) is plotted across trials. Participant bucket  
855 placements were well described by a normative learning model (green line) that adjusts learning rate  
856 according to change-point probability and relative uncertainty, which **(C)** are updated according to  
857 the model on each trial and evolve over time. **D)** Trial-wise measures of neural dissimilarity were  
858 computed on each trial as one minus the correlation coefficient between contiguous trial activations  
859 within a searchlight and regressed onto learning rates from the normative learning model to identify  
860 brain regions with BOLD activations that evolved more rapidly during periods of rapid learning. **E)** A  
861 diverse array of brain regions including occipital regions, dorsomedial prefrontal cortex,  
862 orbitofrontal cortex, and temporal regions displayed neural changes that were positively related to  
863 learning (green clusters). All images are thresholded at  $p = 0.001$  uncorrected.  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881

882 **Figure 2: Representational similarity analysis reveals additional brain regions with**  
883 **representations that evolve more rapidly during periods of learning. A)** In principle, rapid  
884 changes in neural representation coincident with learning might reflect a dynamic state  
885 representation that transitions rapidly at changes in context (see Fig 1a) and evolves more slowly as  
886 subjects develop accurate representations of the context. **B)** This would lead to greater trialwise  
887 dissimilarity immediately after change points in task context (blue line indicates simulated trialwise  
888 dissimilarity, red dashed lines indicate change points), but also to **(C)** unique patterns of dissimilarity  
889 across non-adjacent trials. **D)** A searchlight representational similarity analysis to identify such  
890 patterns revealed a constellation of regions (red) that overlapped substantially with that identified in  
891 the trialwise similarity analysis (orange; conjunction depicted in yellow), and also included  
892 additional regions such as left motor cortex, visual cortex, and hippocampus. All images are  
893 thresholded at  $p = 0.001$  uncorrected.

894  
895

896 **Figure 3: Dissociable explanations for task-driven changes in trialwise dissimilarity. Left:**  
897 Context changes could affect different sorts of representations that are thought to be involved in task  
898 performance. A change in context could elicit a large representational change (arrows) in the  
899 behavioral policy **(A)**, an internal assessment of change-point probability **(B)**, the current level of  
900 relative uncertainty **(C)**, or a latent state that shifts in proportion to learning **(D)**. *Middle:* Each of  
901 these representations would predict increased trialwise dissimilarity after change points (top, red  
902 dotted lines indicate change points). However, dissimilarity matrices constructed across all trials  
903 (adjacent and non-adjacent) reveal unique representational profiles for each source of change-point  
904 related dissimilarity (bottom). *Right:* Patterns of voxel activations across trials revealed an  
905 anatomical dissociation between representations of behavioral policy **(A; left motor cortex)**, change-  
906 point probability **(B; occipital cortex)**, relative uncertainty **(C; widespread)**, and shifting latent states  
907 **(D; orbitofrontal cortex)**. Brain images in each panel reflect t-statistic maps thresholded at  $p < 0.01$   
908 after correction for multiple comparisons. For analogous results using an alternative pre-processing  
909 pipeline (no smoothing before RSA), see extended data figure 3-1.

910  
911  
912  
913  
914  
915  
916  
917

*Peak voxel locations (spatial smoothing before RSA)*

Coefficient	Voxels	Max t-value	X	Y	Z	Label
Behavioral policy	841	6.37	27	-60	-18	Temporal occipital fusiform
	389	6.03	-37	-21	58	Left precentral gyrus (left motor)
Change-point probability	3795	8.13	12	-93	-6	Occipital pole
Uncertainty	29941	11.4	-4	-63	49	Precuneus
	local max	9.4	-22	-90	-15	Occipital fusiform gyrus
	local max	8.6	9	22	37	Anterior cingulate cortex
	local max	8.3	15	-54	1	Lingual gyrus
	local max	8	51	-39	55	Supramarginal gyrus
	local max	8	48	16	1	Insula
Shifting latent state	869	6.02	-61	-24	-24	Inferior temporal gyrus (posterior)
	231	5.48	21	-69	67	Occipitoparietal cortex
	220	5.56	-16	49	-15	Left OFC
	220	5.2	-28	-48	52	Superior parietal lobule
	199	5	27	43	-18	Right OFC
	181	5.6	-13	-93	-9	Occipital pole

918 Table 1: **Peak voxel locations corresponding to behavioral policy, relative uncertainty, change-**  
 919 **point probability and shifting latent state representations.** Cluster size (in voxels), maximum (t-  
 920 statistic) and MNI coordinates for each cluster from the competing computations RSA analysis on  
 921 spatially smoothed data surviving multiple comparisons correction.  
 922  
 923  
 924  
 925  
 926  
 927  
 928  
 929  
 930  
 931  
 932  
 933  
 934  
 935  
 936  
 937

938

*Peak voxel locations (spatial smoothing after RSA)*

Coefficient	Voxels	Max t-value	X	Y	Z	Label
Behavioral policy	1058	6.5	27	-57	-18	R fusiform cortex
	295	5.2	-34	-27	64	L precentral gyrus (motor)
Change-point probability	4191	10.0	21	-90	7	Occipital pole
Uncertainty	29582	10.5	-7	-66	52	Precuneus
	local max	9.0	-19	-84	-18	L Occipital Fusiform
	local max	8.6	-1	-39	58	Postcentral Gyrus
	local max	8.5	30	13	61	R Middle Frontal Gyrus
Shifting latent state	local max	8.4	6	16	52	Paracingulate Gyrus
	3581	5.5	-58	-6	-33	L Middle temporal Gyrus
	2096	6.0	-37	64	-3	Orbitofrontal Cortex
	1290	6.0	-19	-72	64	Sup. Lateral Occ. Complex
	443	4.4	60	-6	-36	R Middle Temporal Gyrus

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

Table 2: **Peak voxel locations corresponding to behavioral policy, relative uncertainty, change-point probability and shifting latent state representations.** Cluster size (in voxels), maximum (t-statistic) and MNI coordinates for each cluster from the competing computations RSA analysis on unsmoothed data surviving multiple comparisons correction (spatial smoothing performed on RSA coefficients before multiple comparisons correction; Fig 3-1 extended data).



*Shifting latent state robustness tests*

Region/Model	Mean Beta	t-value	p-value (uncorrected)
<b>Inferior temporal gyrus (-61, -24, -24)</b>			
Pre-Whitened	0.0375	3.68	8.78e-4
Minimal Model	0.0693	4.57	7.37e-5
Time-Shifted	0.0729	4.19	2.14e-4
<b>Occipitoparietal cortex (21, -69, 67)</b>			
Pre-Whitened	0.0624	3.16	.00347
Minimal Model	0.0372	1.13	.265
Time-Shifted	0.0859	4.09	2.81e-4
<b>Left orbitofrontal cortex (-16, 49, -15)</b>			
Pre-Whitened	0.0256	2.27	.0304
Minimal Model	0.0517	3.43	.00172
Time-Shifted	0.0720	3.98	3.91e-4
<b>Superior parietal lobule (-28, -48, 52)</b>			
Pre-Whitened	0.0175	1.82	.0792
Minimal Model	0.0116	0.547	.588
Time-Shifted	0.0656	4.22	2.00e-4
<b>Right orbitofrontal cortex (27, 43, -18)</b>			
Pre-Whitened	0.0271	2.18	.0367
Minimal Model	0.0586	3.93	4.45e-4
Time-Shifted	0.0640	4.06	3.11e-4
<b>Occipital pole (-13, -93, -9)</b>			
Pre-Whitened	0.0243	2.68	.0116
Minimal Model	0.0539	3.47	.00153
Time-Shifted	0.0426	3.08	.00435

956

957

958 Table 3: Robustness checks in the regions-of-interest that showed a significant effect of  
959 shifting latent state (from peak voxel of clusters reported in table 1). Peak-centered  
960 spheres were re-analyzed in three ways. The "Pre-Whitened" analysis used unsmoothed  
961 voxels that were spatially pre-whitened (Walther et al., 2016). The "Minimal Model"  
962 analysis used a regression model that only contained an intercept, the latent state  
963 predictor, and 15 off-diagonal autocorrelation terms. The "Time-Shifted" analysis used a  
964 time-shifted "shifting latent state" regressor in which representations at the time of  
965 outcome on a given trial were modeled as reflecting the beliefs that would guide  
966 behavior on the subsequent trial. This was offset by one trial from our original analysis,  
967 which assumed that representations upon viewing an outcome would reflect the beliefs  
968 that were formed in anticipation of that outcome, rather than the updated ones that  
969 incorporated it.

970

971

972 Extended Data Figure 3-1: RSA results are robust to spatial smoothing. This extended data figure is  
973 an exact replication of figure 3, except that the analysis differed in the following ways: 1) RSA was  
974 performed on BOLD data that had not been spatially smoothed, and 2) spatial smoothing with a 6mm  
975 FWHM Gaussian kernel was applied to the coefficient maps resulting from RSA. Left: Context changes  
976 could affect different sorts of representations that are thought to be involved in task performance. A  
977 change in context could elicit a large representational change (arrows) in the behavioral policy (A),  
978 an internal assessment of change-point probability (B), the current level of relative uncertainty (C),  
979 or a latent state that shifts in proportion to learning (D). Middle: Each of these representations would  
980 predict increased trialwise dissimilarity after change points (top, red dotted lines indicate change  
981 points). However, dissimilarity matrices constructed across all trials (adjacent and non-adjacent)  
982 reveal unique representational profiles for each source of change-point related dissimilarity  
983 (bottom). *Right:* Patterns of voxel activations across trials revealed an anatomical dissociation  
984 between representations of behavioral policy (A; left motor cortex), change-point probability (B;  
985 occipital cortex), relative uncertainty (C; widespread), and shifting latent states (D; orbitofrontal  
986 cortex). All maps are thresholded at  $p < 0.01$  after correction for multiple comparisons, except the  
987 behavioral policy map in which this threshold was increased to include display of the motor cortical  
988 representation for which the cluster corrected p value was 0.011. Statistics for significant clusters are  
989 reported in table 2.

990  
991  
992





