

1
2 **Dissociable forms of uncertainty-driven representational change across**
3 **the human brain.**
4
5
6
7

8 Matthew R. Nassar¹, Joseph T. McGuire², Harrison Ritz¹ and Joseph Kable³
9

10
11
12
13
14 ¹Department of Cognitive, Linguistic, and Psychological Sciences; Carney
15 Institute for Brain Science, Brown University, Providence RI 02912-1821

16 ²Department of Psychological & Brain Sciences; Boston University, Boston MA
17 02215

18 ³Department of Psychology; University of Pennsylvania, Philadelphia PA 19143
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37

38 Corresponding Author:

39 Matthew R. Nassar

40 Department of Cognitive, Linguistic and Psychological Sciences

41 Brown University

42 Providence, RI 02912-1821

43 Phone: 607-316-4932

44 E-mail: matthew_nassar@brown.edu
45

46

47 **Abstract**

48

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

72

73 **Introduction**

74

75 Neural populations in rodent prefrontal cortex can undergo abrupt changes
76 in firing concomitant with changes in performance in rule-based tasks (Durstewitz
77 et al., 2010; Powell and Redish, 2016). Similar phenomena have been observed in
78 the multi-voxel patterns in human fMRI data preceding changes in task strategy,
79 leading to the notion that such changes might correspond to an "aha moment" at
80 which the brain reorganizes to produce a new task set (Schuck et al., 2015). In
81 rodent learning tasks that involve discontinuously changing reward contingencies,
82 abrupt changes in firing of neurons in medial frontal cortex are observed more
83 frequently during periods of uncertainty, during which animals appear to be
84 searching for the best behavioral policy (Karlsson et al., 2012). It is unclear to what
85 extent such phenomena are specific to medial frontal populations, or to what extent
86 they might have an analog in human learning. Furthermore, while these "network
87 resets" during periods of uncertainty are thought to play a role in behavioral
88 flexibility in changing environments (Tervo et al., 2014) the exact computational
89 role of abrupt changes in such neural representations remains unknown.

90

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

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

96 Alternatively, reset phenomena might result from representation of higher-
97 order computational variables used to appropriately calibrate the rate of learning.
98 Recent work has highlighted a number of computational variables that are
99 important for successful learning in the presence of discontinuous environmental
100 changes (change points). In particular, humans tend to increase rates of learning
101 according to the probability with which a given outcome reflects a change point in
102 the behavioral contingency (*change-point probability*) and according to the relative
103 imprecision of their estimate of the current contingency (*relative uncertainty*)
104 (Nassar et al., 2010; 2012). These computational variables both increase following
105 change-points, albeit with different dynamics, to mediate rapid incorporation of
106 new information during and after periods of environmental change. Change-point
107 probability and relative uncertainty correlate with BOLD responses across a wide
108 swath of brain regions including some that jointly reflect both variables and some
109 that uniquely reflect either change-point probability or uncertainty (McGuire et al.,
110 2014). In principle, neural representations of either computational factor might
111 involve patterns of activation that mimic “network reset” phenomena, yet this
112 possibility has never been tested directly.

113 Another signal that might give rise to reset-like dynamics is a continuously
114 evolving latent state representation. Latent states, which represent the relevant
115 behavioral context in cases where it is not directly observable, can improve learning
116 in the face of abstract stimulus categories or repeated episodes by efficiently
117 partitioning learning across distinct behaviorally relevant contexts (Gershman and
118 Niv, 2010). While previous work has focused primarily on the advantage of such
119 representations for rapid reinstatement of previously learned behaviors (Gershman
120 et al., 2010; Wilson et al., 2014), another advantage of such representations is that
121 they could facilitate rapid disengagement from established behaviors that are no
122 longer relevant. By appropriately partitioning data collected over time in a changing
123 environment, such a mechanism could aid learning even if previously encountered
124 environmental states do not recur. To accomplish this, such a latent state
125 representation would need to evolve faster after a period of environmental change
126 in order to effectively disengage from the previous behavioral context (Prescott
127 Adams and MacKay, 2007; Wilson et al., 2010). While previous work has suggested
128 that orbitofrontal cortex (OFC) might represent latent task states (Wilson et al.,
129 2014; Schuck et al., 2016), it is unclear whether such representations transition
130 dynamically during periods of rapid learning as would be necessary to efficiently
131 mediate disengagement of learned responses that are rendered irrelevant by
132 environmental change.

133 Here we examined whether and where uncertainty-linked network resets are
134 observable in human fMRI data, and evaluated the most likely computational
135 explanation for these phenomena in individual brain regions. We did so using a
136 multistep approach. First, we identified signals that change rapidly from trial to trial
137 during periods of uncertainty and rapid learning and potentially correspond to

138 network resets (Karlsson et al., 2012). Second, we generalized this notion of
139 representational change across pairs of non-consecutive trials using
140 representational similarity analysis (RSA) (Nili et al., 2014). Third, we formalized a
141 set of candidate computational explanations for network-reset phenomena and
142 allowed these explanations to compete to explain multivariate brain activity (Kragel
143 et al., 2018).

144 We observed rapid changes in multivariate activity patterns across
145 widespread cortical regions during periods of uncertainty and rapid learning. Using
146 RSA, we showed that patterns in motor regions were best described as reflecting
147 behavioral policy, patterns of activation in occipital regions were best described as
148 registering the occurrence of change-points, and patterns across much of the rest of
149 the cortex appeared to reflect uncertainty. However, patterns of activation in a small
150 number of regions including OFC were most consistent with dynamic latent state
151 representations, suggesting a possible role for the OFC in translating learning
152 signals into state changes that effectively disengage from behaviors learned in
153 contexts that are no longer relevant.

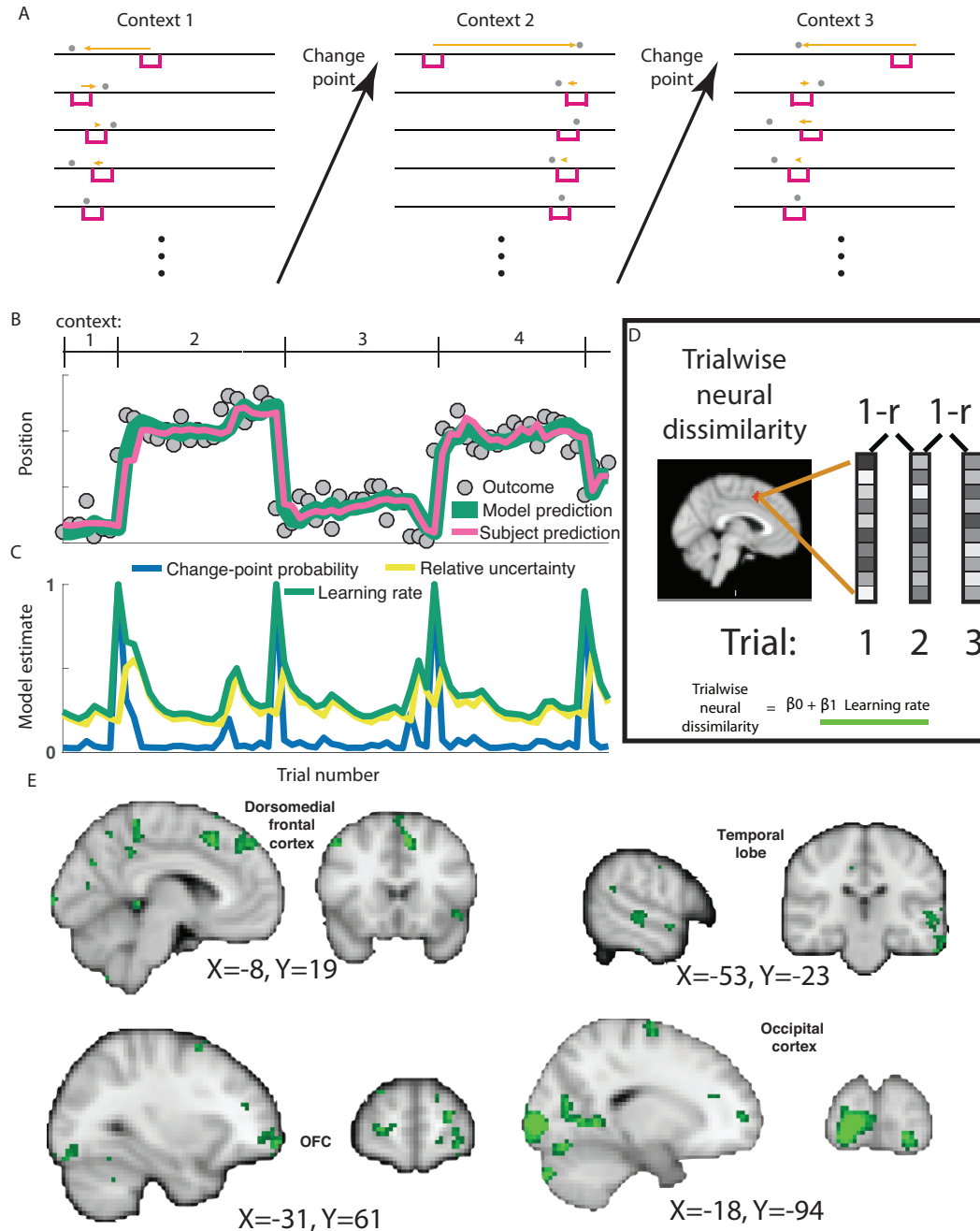
154

155

156 **Results**

157 To examine how neural signals change during periods of uncertainty we re-
158 analyzed data from a previously published study that included recordings of fMRI
159 BOLD signal and behavioral responses of human participants in a predictive
160 inference task (McGuire et al., 2014). Participants played a video game in which they
161 tried to get as many coins as possible (redeemable for money) by catching bags of
162 coins dropped from a hidden helicopter in the sky. Thus, on each task trial,
163 participants estimated the state of an unobservable variable (the position of a
164 helicopter) based on the history of an observable variable (the position of bags
165 dropped from that helicopter) (McGuire et al., 2014). The task included abrupt
166 change points at which the position of the helicopter was resampled from a uniform
167 distribution, which forced participants to rapidly revise beliefs about the helicopter
168 location in order to maintain successful task performance. Here we refer to periods
169 of consistent helicopter position as contexts (Fig 1a), such that the task could be
170 described as requiring dynamic belief updating both within (Fig 1a; vertical) and
171 across (Fig 1a; horizontal) contexts.

172



173
174
175
176
177
178
179
180
181
182
183
184

Figure 1: Trialwise neural dissimilarity is increased after change-points during periods of rapid learning for multiple brain regions. **A)** Participants were asked to move a bucket (pink rectangle) on each trial to the location most likely to deliver a reward (in the form of a bag containing coins). On each trial (stacked vertically) the participant would observe an outcome (bag location; gray circle) that they could use to update their bucket placement for the subsequent trial (orange arrow). Most contiguous trials were generated from the same context, which was defined by a fixed outcome distribution, however at occasional change points, the context (mean outcome location) shifted abruptly and unpredictably. **B)** An example sequence of outcomes (gray circles) and corresponding participant bucket placements (pink line) is plotted across trials. Participant bucket placements were well described by a normative learning model (green line) that adjusts learning rate according to change-point probability and relative uncertainty, which **(C)** are updated according to

185 the model on each trial and evolve over time. **D)** Trial-wise measures of neural dissimilarity were
186 computed on each trial as one minus the correlation coefficient between contiguous trial activations
187 within a searchlight and regressed onto learning rates from the normative learning model to identify
188 brain regions with BOLD activations that evolved more rapidly during periods of rapid learning. **E)** A
189 diverse array of brain regions including occipital regions, dorsomedial prefrontal cortex,
190 orbitofrontal cortex, and temporal regions displayed neural changes that were positively related to
191 learning (green clusters). All images are thresholded at $p = 0.001$ uncorrected.

192

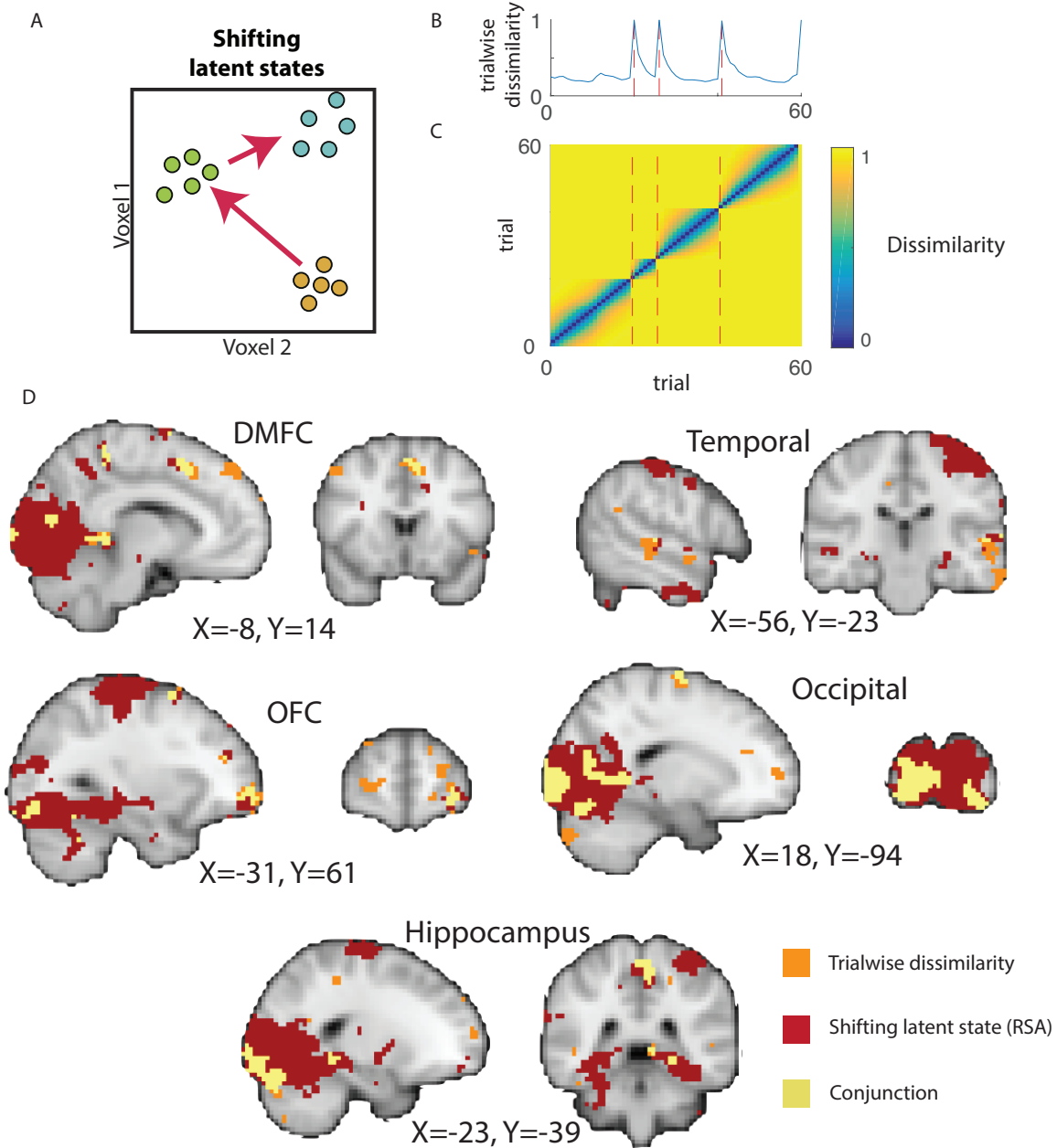
193

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

209 The trial-to-trial dissimilarity in multivariate voxel activation patterns was
210 related to the dynamic learning rates prescribed by the normative model (Fig 1d).
211 Trial wise neural dissimilarity was computed for each pair of sequentially adjacent
212 trials using a whole brain searchlight procedure and regressed onto an explanatory
213 matrix that included model-based estimates of dynamic learning rates. A
214 constellation of regions showed patterns of activation that changed more rapidly
215 during periods of rapid learning after change points (Fig 1e). These regions included
216 OFC, but also clusters in dorsomedial frontal cortex (DMFC), occipital cortex, and the
217 temporal lobe. Thus, with a simple measure of representational change, we
218 identified neural signals whose representations updated more rapidly during
219 periods of learning in multiple brain regions (cf. (Karlsson et al., 2012)).

220 We next exploited representational similarity analysis (RSA) to extend and
221 generalize the analysis above by incorporating information about the pairwise
222 dissimilarity for all pairs of trials, not merely adjacent trial pairs. We hypothesized
223 that the dissimilarity in neural representation for any pair of trials would depend on
224 the cumulative amount of learning expected to occur between them under the
225 normative model (see Methods). The hypothesized pattern of dissimilarity across
226 trials is equivalent to what we would expect from a latent state representation that
227 shifted rapidly at abrupt context transitions and concomitant periods of rapid
228 learning, but remained relatively stable in periods when the statistics of the
229 environment were stationary (Fig 2a). The pattern of dissimilarities predicted
230 across adjacent trials using this strategy is exactly equivalent to the learning rates
231 that served as the explanatory variable in the previous analysis (Fig 2b), but this

232 generalization also makes predictions about the pattern of dissimilarities that would
233 be observed across non-adjacent trials (Fig 2c). We used a searchlight to identify
234 brain regions in which the neural dissimilarity matrix was positively associated with
235 this hypothetical “shifting state representation” hypothesis matrix while controlling
236 for fixed autocorrelation in the similarity structure (see Methods). A significant
237 association was observed in a set of regions that overlapped with the results from
238 the trial-wise dissimilarity analysis, including clusters in OFC, DMFC, occipital, and
239 temporal regions (Fig 2d). As might be expected by the increased power owing to
240 the non-adjacent trial comparisons afforded by RSA analysis, we also identified
241 additional regions that were not clearly indicated by our previous analysis including
242 a number of visual regions, left motor cortex, and bilateral hippocampus (Fig 2d).



243
244
245
246
247
248
249
250
251

Figure 2: Representational similarity analysis reveals additional brain regions with representations that evolve more rapidly during periods of learning. **A)** In principle, rapid changes in neural representation coincident with learning might reflect a dynamic state representation that transitions rapidly at changes in context (see Fig 1a) and evolves more slowly as subjects develop accurate representations of the context. **B)** This would lead to greater trialwise dissimilarity immediately after change points in task context (blue line indicates simulated trialwise

252 dissimilarity, red dashed lines indicate change points), but also to **(C)** unique patterns of dissimilarity
253 across non-adjacent trials. **D)** A searchlight representational similarity analysis to identify such
254 patterns revealed a constellation of regions (red) that overlapped substantially with that identified in
255 the trialwise similarity analysis (orange; conjunction depicted in yellow), and also included
256 additional regions such as left motor cortex, visual cortex, and hippocampus. All images are
257 thresholded at $p = 0.001$ uncorrected.

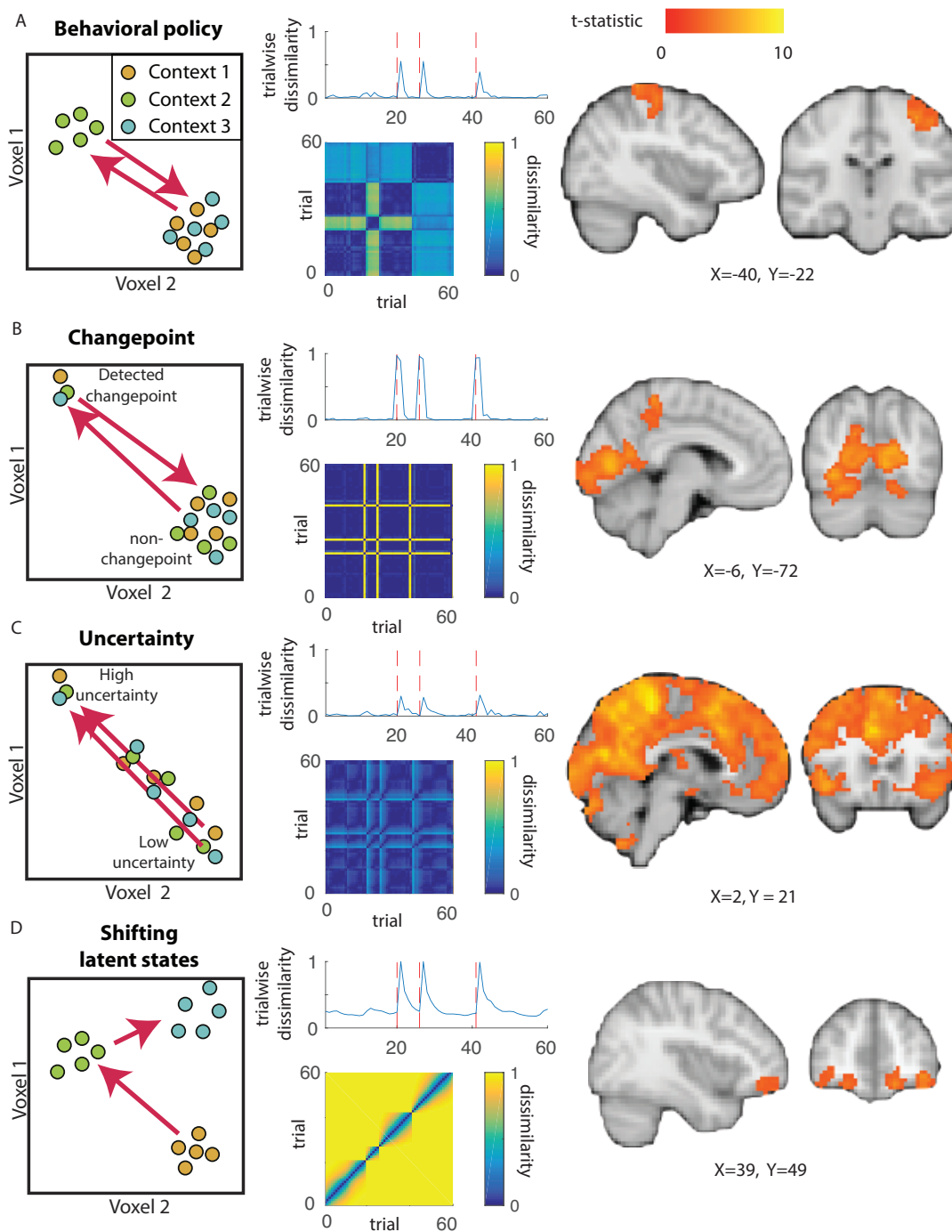
258

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

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

274 A second possible explanation for rapid representational change after change
275 points is that the representations could reflect the current level of change-point
276 probability or relative uncertainty. Change-point probability changes most
277 dramatically at a change in the context (Fig 1c), leading to predicted trialwise neural
278 dissimilarity time courses that do the same (Fig 3b). The level of relative uncertainty
279 changes most rapidly immediately after change-points (Figure 1c), and a neural
280 representation of relative uncertainty should do the same (Fig 3c). However, either
281 of these representations should return to a fixed pattern for all epochs across the
282 experimental session that share the same level of change-point probability or
283 relative uncertainty, irrespective of the current helicopter position (Fig 3b-c).

284



285
286
287
288
289

290 Figure 3: **Dissociable explanations for task-driven changes in trialwise dissimilarity.** *Left:*
291 Context changes could affect different sorts of representations that are thought to be involved in task
292 performance. A change in context could elicit a large representational change (arrows) in the
293 behavioral policy (**A**), an internal assessment of change-point probability (**B**), the current level of
294 relative uncertainty (**C**), or a latent state that shifts in proportion to learning (**D**). *Middle:* Each of
295 these representations would predict increased trialwise dissimilarity after change points (top, red
296 dotted lines indicate change points). However, dissimilarity matrices constructed across all trials
297 (adjacent and non-adjacent) reveal unique representational profiles for each source of change-point
298 related dissimilarity (bottom). *Right:* Patterns of voxel activations across trials revealed an
299 anatomical dissociation between representations of behavioral policy (**A**; left motor cortex), change-
300 point probability (**B**; occipital cortex), relative uncertainty (**C**; widespread), and shifting latent states
301 (**D**; orbitofrontal cortex).

302

303 A final computational explanation for rapid representational changes after
304 change points is that such a signal may reflect a latent state that is used to partition
305 learning across distinct contexts (Wilson et al., 2014). For example, each new
306 helicopter position could be reasonably thought of as a new temporal context,
307 during which learning from prior contexts should be discounted to minimize
308 interference (Fig 1a). Since the helicopter position cannot be resolved exactly, such
309 a context representation would be expected to evolve over time in proportion to the
310 rate of learning about the current context. As described in Figure 2, this would lead
311 to latent state representations that change rapidly at change points and immediately
312 afterwards and change only minimally during periods of prolonged stability (Fig
313 3d). Unlike the other computational factors discussed above, a latent state
314 representation would not necessarily exhibit any systematic similarity relation
315 between one context and another – as our task did not include situations in which
316 the helicopter returned exactly to a previously occupied position. Such a latent state
317 signal might provide an evolving substrate to which outcomes could be linked in
318 order to achieve rational adjustments of learning.

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

329 However, the four candidate representations differed drastically in their
330 predictions about the dissimilarity for non-adjacent pairs of trials. We constructed
331 hypothesis matrices for each candidate representation by considering the expected
332 difference in the computation of interest across all possible pairs of trials. These
333 hypothesis matrices highlight qualitative features of each candidate computation;
334 behavioral policy frequently undergoes abrupt shifts but often takes on a similar
335 value to a previous state, change-point probability highlights differences between
336 change point and non-change point trials, relative uncertainty highlights the
337 differences between high relative uncertainty and other trials, and shifting latent

338 states capture differences largely near the diagonal (Fig 3, middle column, bottom).
339 Consistent with these qualitative differences, correlations between the hypothesis
340 matrices for the different candidate representations were relatively low (all
341 pairwise $r < 0.16$), suggesting that the candidate representations could be efficiently
342 distinguished when considering the entire pairwise dissimilarity matrix.

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

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

363 The relationship between the neural dissimilarity in OFC and the
364 dissimilarity structure predicted by a shifting latent state signal was robust to
365 specific analysis choices. Patterns of activation in right and left OFC clusters were
366 positively related to shifting latent state predictions in the context of our
367 representational similarity regression analysis when using alternative pre-
368 processing strategies such as omitting smoothing (Table 2) or including a spatial
369 pre-whitening procedure (Table 3), both of which emphasize the high frequency
370 components of the spatial pattern (Walther et al., 2016). The observed effects were
371 not driven by relationships between additional explanatory variables included in
372 the regression model, as exclusion of other explanatory variables yielded very
373 similar relationships (Table 4). It is noteworthy that this was not true of all clusters
374 that survived whole-brain correction in our representational similarity regression
375 analysis; clusters identified in left superior parietal lobule and right occipital cortex
376 were not related to the shifting latent state predictions in isolation (Table 4).
377 Furthermore, the relationship between shifting latent state predictions and OFC
378 patterns of activation was also robust to our assumptions about the exact timing of
379 learning; a time shifted version of the shifting latent state hypothesis matrix that
380 assumed learning occurred immediately upon observing a trial outcome could also
381 describe similarity patterns observed in right and left OFC (Table 5).

382 In summary, while we found a number of regions that showed rapidly
383 changing representations during periods of uncertainty following a context change,

384 these reset-like phenomena were due to dissociable computational explanations.
385 While a few regions were implicated in representing behavioral policy or change-
386 point probability, most of these regions reflected relative uncertainty, and a smaller
387 subset of regions including OFC were consistent with representing a latent state that
388 is adjusted according to changes in context.

389

390 **Discussion**

391

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

400 Our analyses revealed distinct explanations for rapid representational
401 changes in different brain regions. Focal representations of behavioral policy and
402 change-point probability were identified in motor and visual cortex respectively,
403 while widespread representations of relative uncertainty were observed throughout
404 the brain. In addition, a small number of brain areas including the OFC had patterns
405 of activation consistent with a form of shifting latent state representation that could
406 speed disengagement from well-learned responses in a changing context.

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

415 Two other computations that we identified using this approach, change-point
416 probability and relative uncertainty, had been the focus of a previous paper using
417 this same dataset (McGuire et al., 2014). In the case of change-point probability,
418 both univariate and RSA analyses revealed occipital cortex and precuneus as the
419 locus of neural representation (see Figure 2c and (McGuire et al., 2014)). However,
420 relative uncertainty representations identified using RSA were considerably more
421 widespread than those identified through univariate activations (see Figure 2c and
422 (McGuire et al., 2014)). This broader set of areas included some regions that were
423 activated in the univariate analysis (e.g., DMFC), some that were deactivated in the
424 univariate analysis (e.g., ventromedial prefrontal cortex), and some that were not
425 identified in univariate analyses at all (e.g., temporal cortex). The near-ubiquitous
426 cortical representation of relative uncertainty revealed by RSA is somewhat
427 analogous to the widespread representations of reward prediction errors that have
428 been identified using multivariate fMRI analysis methods (Vickery et al., 2011).
429 Interestingly, both reward prediction errors and relative uncertainty have been

430 suggested to be signaled through brainstem neuromodulatory systems that could
431 potentially have widespread effects throughout the brain (Schultz, 1997; Yu and
432 Dayan, 2005; Doya, 2008; Nassar et al., 2012).

433 In addition to providing a more sensitive tool to identify well-specified
434 computational variables, RSA also allowed us to look for patterns of activity that
435 could not easily be detected in univariate analyses. In particular, it allowed us to
436 look for neural representations of a dynamically shifting state representation,
437 without making strong assumptions about what the signal would look like at any
438 given moment. It has been proposed that state representations provided by the OFC
439 might serve to hasten learning in environments that include a small number of
440 repeated contexts (Gershman and Niv, 2010; Wilson et al., 2014; Schuck et al.,
441 2016). Here we hypothesized that shifts in the same state representations might
442 implement the rapid learning that should and does follow change-points in outcome
443 contingencies (Prescott Adams and MacKay, 2007; Nassar et al., 2010; Wilson et al.,
444 2010). Such an implementation could make use of existing computational elements
445 to efficiently partition learned associations that pertain to distinct and unrelated
446 contexts, effectively creating the product partitions necessary for optimal inference
447 amid change-points (Prescott Adams and MacKay, 2007).

448 In line with this idea, we identified signals in orbitofrontal cortex consistent
449 with a shifting state signal that changed more rapidly during periods of learning. A
450 neural population that encoded such a signal would be well positioned to transform
451 a direct representation of dynamic learning rate, such as have been identified in
452 cortical regions (Behrens et al., 2007; Krugel et al., 2009; McGuire et al., 2014) and
453 thought to be broadcast through noradrenergic neuromodulation (Yu and Dayan,
454 2005; Nassar et al., 2012; Browning et al., 2015), into a proportional change in
455 associative strength. Using a learning signal to control the rate of contextual shift
456 could enable a simple associative neural network to accomplish the type of adaptive
457 learning that has previously been modeled as a delta-rule update with a varying
458 learning rate. In such a case, increases in apparent learning would be implemented
459 through changes in the substrate for learning, or the active latent state, rather than
460 by adjusting associative strength per se.

461 Representations of latent state that transition dynamically from one context
462 to the next are similar in spirit to the concept of event segmentation in episodic
463 memory (Ezzyat and Davachi, 2010). Segmenting events is useful in that it can allow
464 memories that are embedded within the same event but separated in time to share
465 associations, while memories that may be closer in time but embedded in separate
466 events are maintained separately, preventing interference (Reynolds et al., 2007).
467 One mechanism through which segmentation could be achieved involves dynamic
468 adjustment of the time-constant in slowly fluctuating temporal context signals to
469 effectively “reset” context at event boundaries (Howard and Kahana, 2002; Howard
470 et al., 2010; Manning et al., 2011). Our data suggest a link between this aspect of
471 episodic encoding and the dynamic adjustments of learning that have been observed
472 at context boundaries (Behrens et al., 2007; Nassar et al., 2010; McGuire et al.,
473 2014). However, aspects of our findings also raise questions about the extent of this
474 link. While our results could be interpreted as supporting roles for OFC and
475 temporal lobe in segmenting contexts, we did not observe the same phenomenon in

476 the hippocampus, which is thought to play a key role in event segmentation (Ezzyat
477 and Davachi, 2014; Hsieh et al., 2014; Shapiro, 2014). Instead, we found that
478 representations in hippocampus, like many other brain regions, were best explained
479 as representing uncertainty itself. One potentially relevant detail is that previous
480 contexts were not systematically re-visited in our task, reducing demands for
481 episodic retrieval. An interesting avenue for future work would be to examine how
482 the representations we identified respond when the context abruptly returns to a
483 previously encountered state, such as might require a form of mental time travel for
484 successful performance (Manning et al., 2011).

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

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

500

501

502 **Methods**

503

504

505 *Behavioral task and analysis*

506

507 For details of the behavioral task and data analysis, see our previous report
508 (McGuire et al., 2014). Briefly, 32 human subjects performed a computerized
509 predictive inference task in an MRI scanner while undergoing functional
510 neuroimaging. Each trial required the subject to move a bucket across the horizontal
511 axis of a screen (starting from a "home position" at the right-hand edge, using a
512 joystick controlled by the right hand) to a location that they believed most likely to
513 be underneath a helicopter that was occluded by clouds and thus not directly
514 observable. On each trial, the helicopter would drop a bag that contained either high
515 value or neutral items. Bag locations were normally distributed and centered on the
516 helicopter location (incentivizing bucket placement under the inferred helicopter
517 location). On the majority of trials (90%) the helicopter would remain in the same
518 location as in the previous trial, but occasionally (10%) the helicopter would
519 relocate to a new position along the horizontal axis of the screen (selected randomly
520 and uniformly).

521

522 *MRI data acquisition and preprocessing*

523 T1-weighted MPRAGE structural images (0.9375 X 0.9375 X 1mm voxels, 192
524 X 256 matrix, 160 axial slices, TI=1100ms, TR=1630ms, TE=3.11ms, flip angle=15°),
525 T2*-weighted EPI functional data (3mm isotropic voxels, 64 X 64 matrix, 42 axial
526 slices tilted 30° from the AC-PC plane, TR=2500ms, TE=25ms, flip angle=75°), and
527 fieldmap images (TR=1000ms, TE=2.69 and 5.27ms, flip angle=60°) were acquired
528 on a 3T Siemens Trio with a 32 channel head coil. Functional data were acquired in
529 4 runs, each of which lasted 9 minutes and 25 seconds (226 images).

530 Data were preprocessed using AFNI (Cox, 1996; 2012) and FSL (Jenkinson et
531 al., 2002; Smith et al., 2004; Jenkinson et al., 2012) in the following steps: 1) slice
532 timing correction (AFNI's *3dTshift*), 2) motion correction (FSL's *MCFLIRT*), 3)
533 fieldmap-based geometric undistortion, alignment with structural images, and
534 registration to the MNI template (FSL's *FLIRT* and *FNIRT*), 4) spatial smoothing with
535 a 6mm FWHM Gaussian kernel (FSL's *fslmaths*), 5) outlier attenuation (AFNI's
536 *3dDespike*), and intensity-scaling by a single grand-mean value in each run (FSL's
537 *fslmaths*). The resulting functional time series was deconvolved to estimate trial
538 activations at the time of the bag drop using the least squares-separate method
539 (Mumford et al., 2012) implemented in Matlab.

540

541 *Multivariate fMRI analysis*

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

546 Trial-to-trial dissimilarity scores were extracted by extracting the $i=j-1$
547 diagonal elements from the dissimilarity matrix, which corresponded to the
548 dissimilarity between adjacent trials (see Figure 1d). The dissimilarity scores were
549 regressed onto an explanatory matrix containing an intercept, and dynamic learning
550 rates prescribed by a normative learning model, yielding one coefficient of interest
551 per subject, per searchlight. Dynamic learning rates were estimated as the sum of
552 change-point probability and relative uncertainty minus their product (see Figure
553 1c; (Nassar et al., 2016)). These latent variables were estimated with a parameter-
554 free normative model that took subject prediction errors as an input according to
555 the following set of recursive equations:

556

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

557

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

558

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

559

560

561 where σ_{μ}^2 is the total variance in beliefs about the helicopter location (the generative
562 mean), σ_N^2 is the variance in the distribution of outcomes (bag drops) around that
563 mean, δ_t is the prediction error, and H is the hazard rate and w is the width of the
564 screen. For a full derivation of the model and terms see (Nassar et al., 2010) and for
565 a complete description of the method for estimating latent variables see (Nassar et
566 al., 2016).

567 In general, change-point probability and relative uncertainty were both
568 increased after change-points, albeit with different latencies, leading to learning
569 rates that decay slowly as a function of time within context. Learning rates
570 quantifying sensitivity to information provided on trial j were aligned with the trial-
571 to-trial dissimilarity between trials j and $j+1$. Thus, our analysis targeted patterns of
572 activity whose degree of change between trials j and $j+1$ reflected normative
573 learning predicted to occur from the outcome presented on trial j . The first 3 trials
574 from each block were removed from analysis as they occurred at the onset of fMRI
575 acquisition.

576 Trial-to-trial dissimilarity analysis described above could be thought of as a
577 special case of the general idea that the similarity between each pair of trials might
578 be inversely related to the learning done between them. Because this pattern of
579 similarity is what might be expected to emerge from a representation of the latent
580 task state, which transitions abruptly from one context to the next and remains
581 relatively stable after many trials in a well learned context, we will refer to it as the
582 shifting latent state dissimilarity matrix. The hypothesis matrix for shifting latent
583 states was generated by computing the extent to which the inference on trial i would
584 factor into the inference on trial j , assuming normative learning:
585

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

586 where H is the shifting latent state dissimilarity matrix and α is the learning rate
587 prescribed by a normative model (Nassar et al., 2010), such that more prescribed
588 learning between two trials corresponded to higher values of α , a smaller product
589 term, and thus a greater dissimilarity. The $i=j-1$ diagonal of this matrix is $1-(1-\alpha_i)$, or
590 just α_i and thus equivalent to the vector of trial-to-trial dissimilarities described
591 above. However, the shifting latent state hypothesis matrix also includes
592 information about other elements in the matrix, potentially offering a more
593 powerful construct to ask a similar question. We examined whether this similarity
594 structure was reflected in the neural dissimilarity between trials in each spherical
595 searchlight. The lower triangle of the neural dissimilarity matrix was regressed onto
596 a hypothesis matrix that included an intercept, the shifting latent state hypothesis
597 matrix (lower triangle), and 15 dummy variables designed to remove the influence
598 of autocorrelation on the coefficient of interest. These autocorrelation terms were
599 derived from 15 off-diagonal binary matrices in which a single off diagonal ($i = j-1$; i
600 $= j-2$; $i = j-3 \dots i = j-15$) was set to one. These matrices were constructed to account
601 for any variance in the neural dissimilarity matrices that could be explained by a

602 fixed signal autocorrelation. To be sure that autocorrelation could not affect our
603 analysis of interest, we also set all elements of the shifting latent state similarity
604 matrix that fell outside of this range (trials separated by more than fifteen trials) to
605 the maximum dissimilarity value.

606 To better understand the computations that give rise to rapid changes in
607 neural patterns during periods of learning after a helicopter relocation, we
608 constructed an exhaustive set of hypothesis matrices and conducted a
609 representational similarity analysis in which these representations could compete
610 to explain structure in neural dissimilarity matrices. This analysis required
611 generating hypothesis matrices for various factors that could relate to task
612 uncertainty, learning, or explain nuisance variance in the dissimilarity matrices.
613 Hypothesis matrices were generated for three additional explanatory variables of
614 interest: 1) subject prediction (behavioral policy), 2) relative uncertainty, 3) change-
615 point probability. We also included six additional nuisance variables: 4) the bag
616 drop's location, 5) signed prediction error (ie, the distance between the prediction
617 and the bag drop), 6) high CPP [to account for patterns of activity that may
618 asymmetrically encode CPP], 7) high RU [to account for patterns of activity that may
619 asymmetrically encode RU], 8) outcome reward value, and 9) task block. For factors
620 1-5 and 8, element (i,j) of the hypothesis matrix corresponded to the absolute
621 difference in that factor on trials i and j. For factor 9, dissimilarity values were set to
622 0 for trials in the same block and 1 for trials in different blocks. Dissimilarity
623 matrices for factors 6 & 7 were computed as one minus the multiplicative
624 interaction of the model variable (6=change-point probability, 7=relative
625 uncertainty) on trials i and j, such that similarity was only hypothesized when the
626 model-derived term took on a high value on both trials. These terms allowed the
627 model to capture asymmetric representations of the two factors governing learning
628 in our model, such as a representation that converged for values of high relative
629 uncertainty but did not show any consistent pattern of activation when relative
630 uncertainty was low.

631 The lower triangle of the neural dissimilarity matrix was extracted and
632 regressed onto an explanatory matrix consisting of an intercept and the lower
633 triangle of all hypothesis/nuisance matrices (including the shifting latent state and
634 nuisance autocorrelation terms), yielding one coefficient per variable, per subject,
635 per searchlight (Chikazoe et al., 2014; Kragel et al., 2018). Group level analyses were
636 conducted by computing t-statistics across subjects for each variable and
637 searchlight. Cluster-based permutation testing using cluster mass with a cluster
638 forming threshold of $p < 0.001$ and an alpha of 0.01 was used to identify significant
639 activations (Nichols and Holmes, 2002).

640
641
642
643
644
645
646
647

648

649 Reference:

650

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

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

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

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

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

661 Doya K (2008) Modulators of decision making. *Nature Neuroscience* 11:410–416.

662 Durstewitz D, Vittoz NM, Floresco SB, Seamans JK (2010) Abrupt Transitions
663 between Prefrontal Neural Ensemble States Accompany Behavioral Transitions
664 during Rule Learning. *Neuron* 66:438–448.

665 Ezzyat Y, Davachi L (2010) What Constitutes an Episode in Episodic Memory?
666 *Psychol Sci* 22:243–252.

667 Ezzyat Y, Davachi L (2014) Similarity Breeds Proximity: Pattern Similarity within
668 and across Contexts Is Related to Later Mnemonic Judgments of Temporal
669 Proximity. *Neuron* 81:1179–1189.

670 Gershman SJ, Blei DM, Niv Y (2010) Context, learning, and extinction. *Psychological*
671 *Review* 117:197–209.

672 Gershman SJ, Niv Y (2010) Learning latent structure: carving nature at its joints.
673 *Current Opinion in Neurobiology* 20:251–256.

674 Howard MW, Kahana MJ (2002) A Distributed Representation of Temporal Context.
675 *Journal of Mathematical Psychology* 46:269–299.

676 Howard MW, Shankar KH, Jagadisan UKK (2010) Constructing Semantic
677 Representations From a Gradually Changing Representation of Temporal
678 Context. *Top Cogn Sci* 3:48–73.

679 Hsieh L-T, Gruber MJ, Jenkins LJ, Ranganath C (2014) Hippocampal Activity Patterns
680 Carry Information about Objects in Temporal Context. *Neuron* 81:1165–1178.

681 Jenkinson M, Bannister P, Brady M, Smith S (2002) Improved optimization for the

- 682 robust and accurate linear registration and motion correction of brain images.
683 *NeuroImage* 17:825–841.
- 684 Jenkinson M, Beckmann CF, Behrens TEJ, Woolrich MW, Smith SM (2012) FSL.
685 *NeuroImage* 62:782–790.
- 686 Karlsson MP, Tervo DGR, Karpova AY (2012) Network resets in medial prefrontal
687 cortex mark the onset of behavioral uncertainty. *Science* 338:135–139.
- 688 Kragel PA, Kano M, Van Oudenhove L, Ly HG, Dupont P, Rubio A, Delon-Martin C,
689 Bonaz BL, Manuck SB, Gianaros PJ, Ceko M, Reynolds Losin EA, Woo C-W,
690 Nichols TE, Wager TD (2018) Generalizable representations of pain, cognitive
691 control, and negative emotion in medial frontal cortex. Nature Publishing Group.
- 692 Krugel LK, Biele G, Mohr PNC, Li S-C, Heekeren HR (2009) Genetic variation in
693 dopaminergic neuromodulation influences the ability to rapidly and flexibly
694 adapt decisions. *Proceedings of the National Academy of Sciences* 106:17951–
695 17956.
- 696 Manning JR, Polyn SM, Baltuch GH, Litt B, Kahana MJ (2011) Oscillatory patterns in
697 temporal lobe reveal context reinstatement during memory search. *Proceedings*
698 *of the National Academy of Sciences* 108:12893–12897.
- 699 McGuire JT, Nassar MR, Gold JI, Kable JW (2014) Functionally dissociable influences
700 on learning rate in a dynamic environment. *Neuron* 84:870–881.
- 701 Mumford JA, Turner BO, Ashby FG, Poldrack RA (2012) Deconvolving BOLD
702 activation in event-related designs for multivoxel pattern classification analyses.
703 *NeuroImage* 59:2636–2643.
- 704 Nassar MR, Bruckner R, Gold JI, Li S-C, Heekeren HR, Eppinger B (2016) Age
705 differences in learning emerge from an insufficient representation of
706 uncertainty in older adults. *Nature Communications* 7:11609.
- 707 Nassar MR, Rumsey KM, Wilson RC, Parikh K, Heasley B, Gold JI (2012) Rational
708 regulation of learning dynamics by pupil-linked arousal systems. *Nature*
709 *Neuroscience* 15:1040–1046.
- 710 Nassar MR, Wilson RC, Heasley B, Gold JI (2010) An approximately Bayesian delta-
711 rule model explains the dynamics of belief updating in a changing environment.
712 *Journal of Neuroscience* 30:12366–12378.
- 713 Nichols TE, Holmes AP (2002) Nonparametric permutation tests for functional
714 neuroimaging: a primer with examples. *Hum Brain Mapp* 15:1–25.
- 715 Nili H, Wingfield C, Walther A, Su L, Marslen-Wilson W, Kriegeskorte N (2014) A
716 Toolbox for Representational Similarity Analysis Prlic A, ed. *PLoS Comput Biol*

- 717 10:e1003553.
- 718 Powell NJ, Redish AD (2016) Representational changes of latent strategies in rat
719 medial prefrontal cortex precede changes in behaviour. *Nature Communications*
720 7:12830.
- 721 Prescott Adams R, MacKay DJC (2007) Bayesian Online Changepoint Detection.
722 eprint arXiv:07103742:-.
- 723 Reynolds JR, Zacks JM, Braver TS (2007) A computational model of event
724 segmentation from perceptual prediction. *Cogn Sci* 31:613–643.
- 725 Schuck NW, Cai MB, Wilson RC, Niv Y (2016) Human Orbitofrontal Cortex
726 Represents a Cognitive Map of State Space. *Neuron* 91:1402–1412.
- 727 Schuck NW, Gaschler R, Wenke D, Heinzle J, Frensch PA, Haynes J-D, Reverberi C
728 (2015) Medial Prefrontal Cortex Predicts Internally Driven Strategy Shifts.
729 *Neuron* 86:331–340.
- 730 Schultz W (1997) A Neural Substrate of Prediction and Reward. *Science* 275:1593–
731 1599.
- 732 Shapiro ML (2014) Time and Again. *Neuron* 81:964–966.
- 733 Smith SM, Jenkinson M, Woolrich MW, Beckmann CF, Behrens TEJ, Johansen-Berg H,
734 Bannister PR, De Luca M, Drobnjak I, Flitney DE, Niazy RK, Saunders J, Vickers J,
735 Zhang Y, De Stefano N, Brady JM, Matthews PM (2004) Advances in functional
736 and structural MR image analysis and implementation as FSL. *NeuroImage* 23
737 Suppl 1:S208–S219.
- 738 Tervo DGR, Proskurin M, Manakov M, Kabra M, Vollmer A, Branson K, Karpova AY
739 (2014) Behavioral Variability through Stochastic Choice and Its Gating by
740 Anterior Cingulate Cortex. *Cell* 159:21–32.
- 741 Vickery TJ, Chun MM, Lee D (2011) Ubiquity and Specificity of Reinforcement
742 Signals throughout the Human Brain. *Neuron* 72:166–177.
- 743 Walther A, Nili H, Ejaz N, Alink A, Kriegeskorte N, Diedrichsen J (2016) Reliability of
744 dissimilarity measures for multi-voxel pattern analysis. *NeuroImage* 137:188–
745 200.
- 746 Wilson RC, Nassar MR, Gold JI (2010) Bayesian online learning of the hazard rate in
747 change-point problems. *Neural Comput* 22:2452–2476.
- 748 Wilson RC, Takahashi YK, Schoenbaum G, Niv Y (2014) Orbitofrontal cortex as a
749 cognitive map of task space. *Neuron* 81:267–279.
- 750 Yu AJ, Dayan P (2005) Uncertainty, neuromodulation, and attention. *Neuron*

751 46:681–692.

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803
804
805
806
807

coefficient	Voxels	Max t	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	Precuneous
	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

Table 1: 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 surviving multiple comparisons correction.

Latent state analysis with unsmoothed voxels

Region	Mean Beta	t-value	p-value (uncorrected)
Left inferior temporal gyrus	0.0663	4.42	1.11e-4
Left superior parietal lobule	0.0491	3.51	.00138
Right occipital cortex	0.1104	4.80	3.76e-5
Left orbitofrontal cortex	0.0541	3.36	.00210
Right orbitofrontal cortex	0.0649	4.08	2.89e-4
Left occipital pole	0.0442	2.97	.00574

Table 2: Regions-of-interest that showed a significant effect of shifting latent state, re-analyzed with unsmoothed voxels.

Latent state analysis with unsmoothed, pre-whitened voxels

Region	Mean Beta	t-value	p-value (uncorrected)
Left inferior temporal gyrus	0.0375	3.68	8.78e-4
Left superior parietal lobule	0.0175	1.82	.0792
Right occipital cortex	0.0624	3.16	.00347
Left orbitofrontal cortex	0.0256	2.27	.0304
Right orbitofrontal cortex	0.0271	2.18	.0367
Left occipital pole	0.0243	2.68	.0116

Table 3: Regions-of-interest that showed a significant effect of shifting latent state, re-analyzed with unsmoothed voxels that were spatial pre-whitened (Walther et al., 2016).

Minimal latent state analysis with unsmoothed voxels

Region	Mean Beta	t-value	p-value (uncorrected)
Left inferior temporal gyrus	0.0693	4.57	7.37e-5
Left superior parietal lobule	0.0116	0.547	.588
Right occipital cortex	0.0372	1.13	.265
Left orbitofrontal cortex	0.0517	3.43	.00172
Right orbitofrontal cortex	0.0586	3.93	4.45e-4
Left occipital pole	0.0539	3.47	.00153

Table 4: Latent state effect in ROIs sensitive to latent state, re-analyzed with unsmoothed voxels and a model that only contained an intercept, the latent state predictor, and 15 off-diagonal autocorrelation terms.

Time shifted latent state analysis

Region	Mean Beta	t-value	p-value (uncorrected)
Left inferior temporal gyrus	0.0729	4.19	2.14e-4
Left superior parietal lobule	0.0656	4.22	2.00e-4
Right occipital cortex	0.0859	4.09	2.81e-4
Left orbitofrontal cortex	0.0720	3.98	3.91e-4
Right orbitofrontal cortex	0.0640	4.06	3.11e-4
Left occipital pole	0.0426	3.08	.00435

Table 5: Shifting latent state effect in ROIs sensitive to shifting latent state, re-analyzed using a time-shifted “shifting latent state” regressor in which representations at the time of outcome on a given trial are modeled as reflecting the beliefs that will guide behavior on the subsequent trial. This is offset by one trial from our original analysis, which assumed that representations upon viewing an outcome would reflect the beliefs that were formed in anticipation of that outcome, rather than the updated ones that incorporated it.