

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

**Title:** Ventromedial prefrontal cortex compression during concept learning

**Brief title:** Ventromedial PFC compression during learning

**Authors:** Michael L. Mack<sup>a</sup>, Alison R. Preston<sup>b,c,d\*</sup>, Bradley C. Love<sup>e,f\*</sup>

**Affiliations:**

<sup>a</sup>Department of Psychology, University of Toronto, Toronto, ON, CA

<sup>b</sup>Department of Psychology, The University of Texas at Austin, Austin, TX, USA

<sup>c</sup>Center for Learning and Memory, The University of Texas at Austin, Austin, TX, USA

<sup>d</sup>Department of Neuroscience, The University of Texas at Austin, Austin, TX, USA

<sup>e</sup>Experimental Psychology, University College London, London, UK

<sup>f</sup>Alan Turing Institute, London, UK

\*Authors contributed equally

**Corresponding Author:**

Michael L. Mack

Department of Psychology

University of Toronto

100 St. George Street, 4<sup>th</sup> Floor

Toronto, Ontario M5S 3G3

mack.michael@gmail.com

**Keywords:** prefrontal cortex; fMRI; attention; category learning; computational modeling

## Abstract

35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48

Prefrontal cortex (PFC) is thought to support the ability to focus on goal-relevant information by filtering out irrelevant information, a process akin to dimensionality reduction. Here, we test this dimensionality reduction hypothesis by combining a data-driven approach to characterizing the complexity of neural representation with a theoretically-supported computational model of learning. We find strong evidence of goal-directed dimensionality reduction within human ventromedial PFC during learning. Importantly, by using model predictions of each participant's attentional strategies during learning, we find that that the degree of neural compression predicts an individual's ability to selectively attend to concept-specific information. These findings suggest a domain-general mechanism of learning through compression in ventromedial PFC.

## 49 Introduction

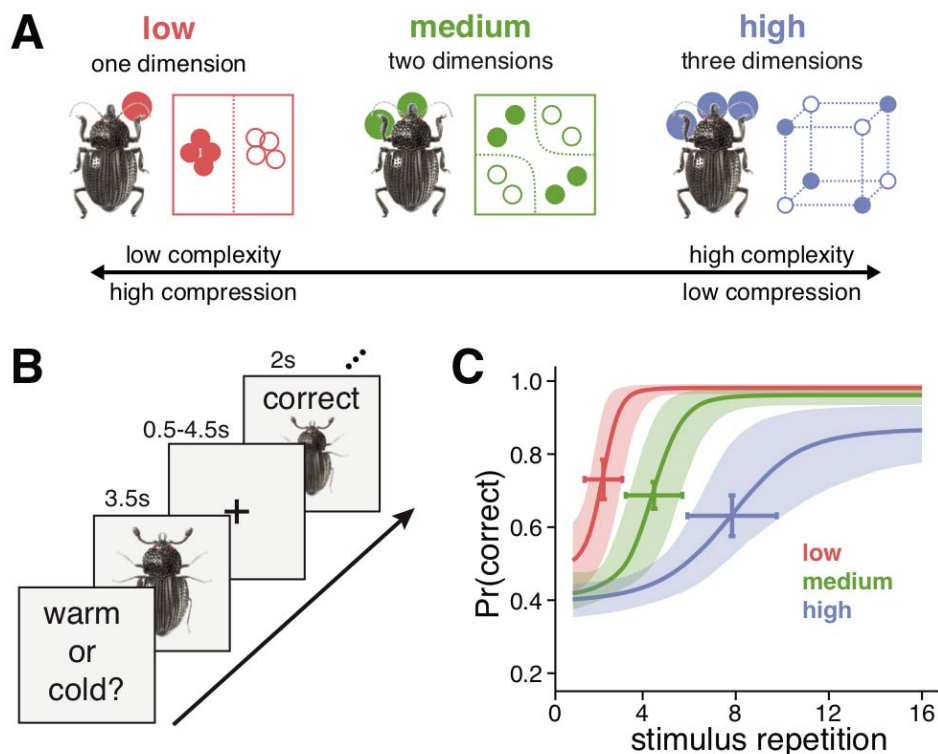
50

51 Prefrontal cortex (PFC) is sensitive to the complexity of incoming information<sup>1</sup> and  
52 theoretical perspectives suggest that a core function of PFC is to focus representation  
53 on goal-relevant features by filtering out irrelevant content<sup>2,3</sup>. In particular, ventromedial  
54 PFC (vmPFC) is thought to represent the latent structures of experience<sup>4,5</sup>, coding for  
55 causal links<sup>6</sup> and task-related cognitive maps<sup>7</sup>. At the heart of these accounts is the  
56 hypothesis that during learning, mPFC may perform data reduction on incoming  
57 information, compressing task-irrelevant features and emphasizing goal-relevant  
58 information structures. This compression process is goal-directed and akin to how  
59 attention in category learning models dynamically selects features that have proven  
60 predictive across recent learning trials<sup>8,9</sup>. Although emerging evidence suggests  
61 structured representations occur in the rodent homologue of vmPFC<sup>10,11</sup>, such coding in  
62 human vmPFC remains poorly understood. Here, we directly assess the data reduction  
63 hypothesis by leveraging an information-theoretic approach in human neuroimaging to  
64 measure how goal-driven learning is supported by attention updating processes in  
65 vmPFC.

66

67 We focused on concept learning, given the recent findings that vmPFC represents  
68 conceptual information in an organized fashion<sup>12,13</sup>. Participants learned to classify the  
69 same insect images (Figure 1A), composed of three features that could take on two  
70 values (thick/thin legs, thick/thin antennae, pincer/shovel mandible), across three  
71 different learning problems<sup>14</sup>. These learning problems were defined by rules that  
72 required consideration of different numbers of features to successfully classify (see  
73 Table 1): the low category complexity problem was unidimensional (e.g., insects living in  
74 warm climates have thick legs, cold climate insects have thin legs), the medium  
75 category complexity problem depended on two features (e.g., insects from rural  
76 environments have thick antennae and shovel mandible or thin antennae and pincer  
77 mandible, urban insects have thick antennae and pincer mandible or thin antennae and  
78 shovel mandible), and the high category complexity problem required all three features  
79 (i.e., each insect's class was uniquely defined by a combination of features). By using  
80 the same stimuli for all three problems, the manipulation of conceptual complexity  
81 allowed us to target goal-specific learning processes.

82



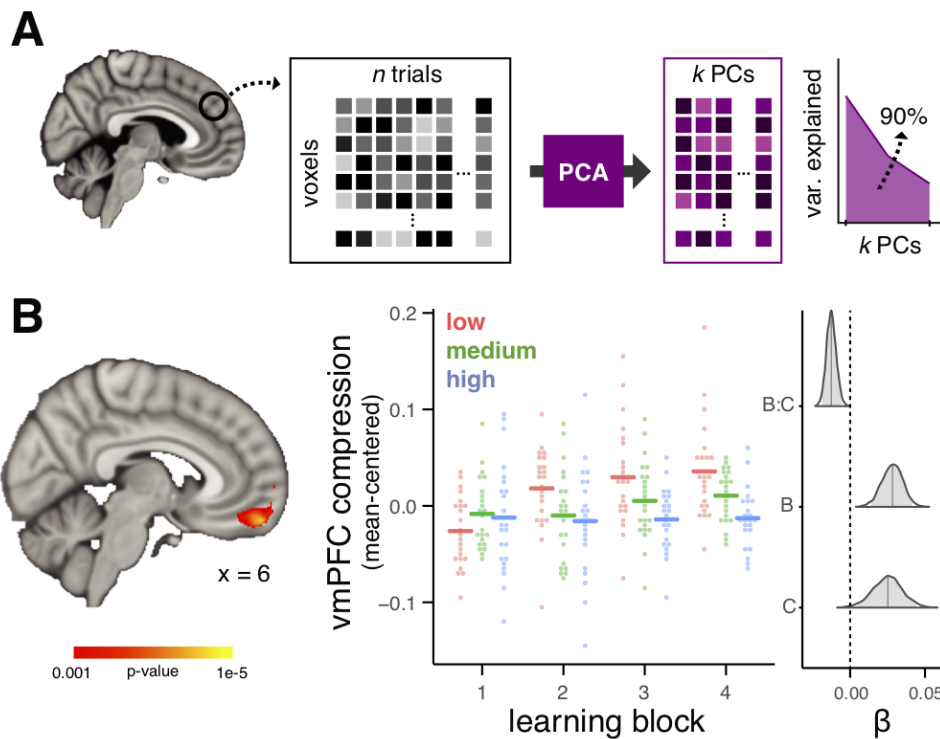
83  
84 **Figure 1:** Experimental schematic and behavioral results (N=23). **A)** The learning problems differed in  
85 rule complexity (see Table S1). The low complexity problem was unidimensional (e.g., antennae size),  
86 medium complexity required a conjunction of two features (e.g., leg size and mandible shape), and high  
87 complexity required all three features. **B)** Learning trials consisted of presentation of a stimulus for 3.5s,  
88 followed by a fixation cross for 0.5-4.5s, and then a feedback display for 2s that included the stimulus,  
89 accuracy of the response, and the correct category. Learning trials were separated by a delay of 2-6s of  
90 fixation. **C)** The probability of a correct response increased across stimulus repetitions. The rate of  
91 learning differed according to the complexity of the problems.

92  
93 This design allows us to directly test whether compression of neural representations  
94 corresponds with the complexity of the problem-specific conceptual structure during  
95 learning. Complexity and compression have an inverse relationship; the lower the  
96 complexity of a conceptual space, the higher the degree of compression. For instance,  
97 in learning the unidimensional problem, variance along the two irrelevant feature  
98 dimensions can be compressed resulting in a lower complexity conceptual space. In  
99 contrast, learning the high complexity problem would result in less compression  
100 because all three feature dimensions must be represented, resulting in a more complex  
101 conceptual space relative to the unidimensional problem. Differences in complexity  
102 across the three learning problems thus provide a means for testing how learning  
103 shapes the dimensionality of neural concept representations. Namely, brain regions  
104 involved in goal-directed data compression should learn by building internal models that  
105 adapt to the complexity of the problems in order to represent information relevant to the  
106 task at hand.

107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126

## Results

To test this prediction, we recorded functional magnetic resonance imaging (fMRI) data while participants learned the three problems and measured the degree that multivoxel activation patterns were compressed through learning using principal component analysis (PCA; Figure 2A), a method for low-rank approximation of multidimensional data<sup>15</sup>. Specifically, trial-level whole brain activation patterns for each insect image were estimated using the LS-S approach<sup>16</sup>. These trial-specific activation patterns were then submitted to PCA and the number of principal components (PC) that were necessary to explain 90% of the variance across trials within a learning block was used to calculate an index of neural compression (i.e., fewer PCs reflects more neural compression). This measure of neural compression was calculated across the whole brain with searchlight methods<sup>17</sup> for each learning block in each problem. We then identified brain regions that reduce dimensionality with learning (i.e., learn to represent the less complex problems with fewer dimensions) by conducting a voxel-wise linear mixed effects regression on the searchlight compression maps. Specifically, at each voxel, we assessed how neural compression changed as a function of learning block and problem complexity and their interaction.



127  
128  
129  
130

**Figure 2:** Neural compression analysis schematic and results (N=23). **A**) Principal component analysis (PCA) was performed on neural patterns evoked for each of  $n$  trial within a learning block. The number of

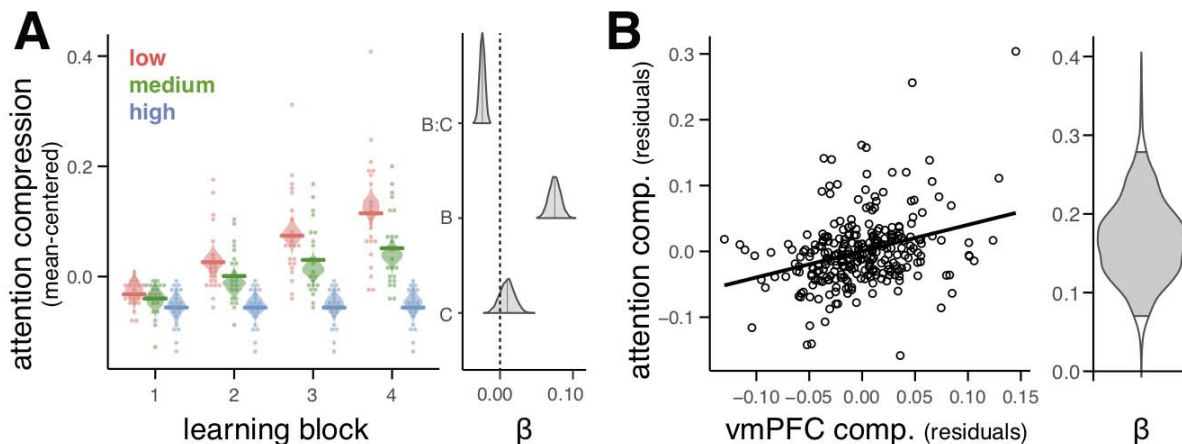
131 principal components (PC) required to explain 90% of the variance ( $k$ ) was used to calculate a neural  
132 compression score ( $1-k/n$ ). We quantified neural compression as a function of problem complexity and  
133 learning block; the interaction of these factors reflects changes in the complexity of neural representations  
134 that emerge with learning. **B)** A whole brain voxel-wise linear mixed effects regression revealed a vmPFC  
135 region that showed a significant interaction between learning block and problem complexity. See  
136 Supplementary Figure 1 for main effect maps of learning block and problem complexity. The nature of the  
137 interaction in the vmPFC region is depicted in the middle panel; points represent compression at the  
138 cluster's peak voxel for each participant and the horizontal lines depict the group average. The right graph  
139 plots the results of a Bayesian linear mixed effects regression of neural compression from the peak voxel  
140 of the vmPFC cluster. Posterior distributions of coefficients from the regression model are depicted for the  
141 factors of learning block (B), complexity (C), and their interaction (B:C). Shaded regions within the  
142 distributions represent 95% high-density intervals. These data and regression results are displayed only  
143 to demonstrate the nature of the interaction effect in the vmPFC cluster and do not represent an  
144 independent statistical analysis.

145  
146 Throughout the entire brain, only a region within vmPFC showed an interaction of  
147 problem complexity and learning block (peak coordinates [4, 54, -18]; 653 voxels; voxel-  
148 wise threshold = 0.001, cluster extent threshold = 0.05; Figure 2B; see Supplementary  
149 Figure 1 and Supplementary Table 1 for main effects of problem complexity and  
150 learning block). The nature of the interaction within this cluster showed that vmPFC  
151 compression corresponded with problem complexity and emerged over learning blocks  
152 (peak:  $\beta_{mean} = -0.013$ , 95% HDI = [-0.019, -0.006],  $P < 0.001$ ). Importantly, the  
153 interaction effect was independent of problem order, individual differences in learning  
154 performance, and remained when looking at only the low and medium complexity  
155 problems (see Methods for details about the voxel-wise regression modeling and control  
156 analyses). Moreover, there was no interaction of learning block and problem complexity  
157 when examining univariate vmPFC activation ( $\beta_{mean} = -1.309$ , 95% HDI = [-82.7, 126.7],  
158  $P = 0.93$ ), ruling out an explanation based on problem difficulty impacting overall neural  
159 activation. Because the stimuli were identical across the three problems, this finding  
160 demonstrates that learning-related compression is goal-specific, with vmPFC requiring  
161 fewer dimensions for less complex goals.

162  
163 To evaluate whether vmPFC compression tracked changes in attentional allocation, we  
164 characterized the participant-specific attentional weights given to each stimulus feature  
165 across the three problems using a computational learning model<sup>8</sup>. Attention weight  
166 compression indexed changes in attentional allocation; low attention compression  
167 indicates equivalent weighting to all three features, whereas high attention compression  
168 indicates attention directed to only one feature. We found that similar to vmPFC  
169 compression, attention compression varied with the interaction of learning block and  
170 conceptual complexity ( $\beta_{mean} = -0.028$ , 95% HDI = [-0.035, -0.020],  $P < 0.001$ ),  
171 suggesting a link between the behavioral and neural signatures of dimensionality  
172 reduction.

173

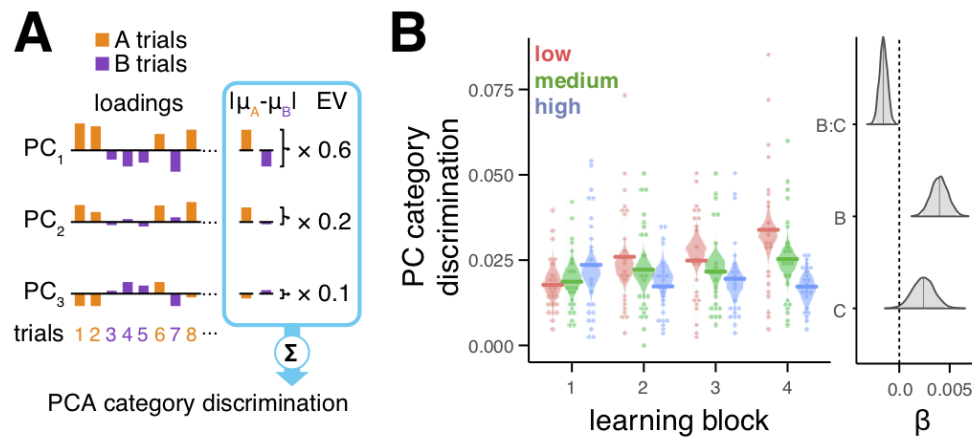
174 To assess this relationship, we evaluated whether the compression of participants'  
175 attention weights was predicted by vmPFC neural compression at the individual  
176 participant level. Specifically, if the ability to compress neural representations in a  
177 problem-appropriate fashion is related to participants' ability to attend to problem-  
178 relevant features, the prediction follows that participants with more neural compression  
179 for a given problem will also show more selective attention, thus higher attention  
180 compression values. This hypothesis was confirmed with a Bayesian mixed effects  
181 regression analysis that took into account differences across learning block, problem  
182 complexity, accuracy, and learning order ( $\beta_{mean} = 0.168$ , 95% HDI = [0.072, 0.277],  $P =$   
183 0.0005; see Figure 3B). Importantly, this relationship remained when restricting analysis  
184 to the low and medium complexity problems, which were counterbalanced for learning  
185 order ( $\beta_{mean} = 0.210$ , 95% HDI = [-0.082, 0.368],  $P = 0.0005$ ). Thus, even after  
186 controlling for differences in neural and attention compression due to learning block and  
187 problem complexity, the degree of problem-specific compression in vmPFC  
188 representations significantly predicted participants' attentional strategies throughout  
189 learning.  
190



191  
192  
193 **Figure 3:** Relationship between model-based attention weighting and vmPFC compression (N=23). **A)**  
194 Attention compression (i.e., degree of attention selectivity) emerged according to feature relevancy  
195 across the problems with highest compression for low complexity followed by medium and high  
196 complexity by the end of learning. Points depict individual participants, horizontal lines are group  
197 averages, and violin plots depict the distribution of posterior predictions from the Bayesian linear  
198 regression. The right panel depicts posterior distributions of regression coefficients for learning block (B),  
199 complexity (C), and their interaction (B:C). The shaded region within the distributions marks the 95% high-  
200 density interval. **B)** vmPFC compression predicted the degree of problem-specific attention weighting  
201 (indexed as attention compression) across learning blocks and problems. Both vmPFC compression and  
202 attention compression are plotted as partial residuals of separate regression models that regress out  
203 factors of learning block, problem complexity, accuracy, and learning order. The solid line depicts the  
204 best-fitting regression line of the partial residuals. The violin plot depicts the posterior distribution of  
205 regression coefficient relating neural compression and attention compression. The shaded area bounded  
206 by black lines within the distribution mark the 95% high-density interval.

207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225

Although vmPFC compression tracks model-based predictions of learning, this link between learning problem-specific coding and neural representation is ultimately indirect. The neural compression findings may be due to learning-related changes in neural representation that highlight within-category similarities and differentiate between-category differences or due to other factors unrelated to the category structure of the problem at hand. To directly assess the degree of category coding present in vmPFC compression, we analyzed how trials loaded onto the PCs. Specifically, we hypothesized that if neural compression is driven by category-specific coding in activation patterns, trials will load on the PCs similarly within category, but differently between category (Figure 4A). In contrast, if neural compression is due to factors not related to category representation, trials from both categories will load similarly on the PCs. We found support for the former (Figure 4B and Table 3) such that category discrimination in PCA loadings increased over learning blocks ( $\beta_{mean} = 0.008$ , 95% HDI = [0.004, 0.011],  $P < 0.001$ ) and was highest for the low followed by medium and high complexity problems ( $\beta_{mean} = -0.003$ , 95% HDI = [-0.0048, -0.0013],  $P < 0.001$ ). Thus, by directly assessing the structure of the PCA results, we find that vmPFC compression is driven by activation patterns that discriminate categories based on current task goals.



226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238

**Figure 4:** Category discrimination in neural compression (N=23). **A)** Category discrimination in PCA neural compression was indexed by the difference in loadings for the two categories. The absolute difference between the average loadings for category A trials (orange) and B trials (purple) for each PC were weighted by the PC's explained variance ratio and summed. **B)** Category discrimination increased across learning blocks with the highest discrimination for low followed by medium and high problem complexity. Points depict individual participants, horizontal lines are group averages, and violin plots depict the distribution of posterior predictions from the Bayesian linear regression. The right panel depicts posterior distributions of regression coefficients for learning block (B), complexity (C), and their interaction (B:C). The shaded region within the distributions marks the 95% high-density interval.

## Discussion



239  
240 By focusing on a mechanism by which vmPFC may form and represent concepts  
241 through goal-sensitive dimensionality reduction, we show that neural representations in  
242 a vmPFC subregion are shaped by experience. And, this shaping is adaptive, promoting  
243 efficient representation of information that focuses on encoding features that are most  
244 predictive of positive outcomes for a given goal. Importantly, by evaluating behavior  
245 through the lens of a theoretically-oriented computational model, we demonstrate that  
246 compression in vmPFC unfolds over the course of learning in a manner consistent with  
247 the learning mechanisms of SUSTAIN<sup>8,9</sup>. These findings provide a quantitative account  
248 of vmPFC's potential role in the coding of efficient schematic models or cognitive  
249 maps<sup>7,12,18,19</sup>, specifically in the conceptual domain.

250  
251 Successfully learning new concepts requires attending to goal-diagnostic features and  
252 ignoring irrelevant information to build abstract representations that capture the  
253 structure defining a concept<sup>8</sup>. Viewed in these terms, concept learning has many  
254 parallels to schema formation, a vmPFC-related function born out of lesion studies in  
255 the memory literature<sup>20</sup>. Schemas are defined as structured memory networks that  
256 represent associative relationships among prior experiences and provide predictions for  
257 new experiences<sup>4,5,21,22</sup>. Schema-related memory behaviors are significantly impacted  
258 by vmPFC lesions. For example, vmPFC lesion patients exhibit a reduced influence of  
259 prior knowledge during recognition of items presented in schematically congruent  
260 contexts compared with healthy controls<sup>23</sup>. Moreover, vmPFC lesions have been  
261 associated with a marked inability to differentiate schema-related concepts from  
262 concepts inappropriate for a given schema<sup>24</sup>. From this work, it is clear that vmPFC is  
263 necessary for retrieving generalized representations built from prior events that are  
264 relevant to current experience. Such guided retrieval of relevant learned representations  
265 is key to building new concepts.

266  
267 A key proposal of the SUSTAIN computational model we leveraged is that concept  
268 learning is decidedly goal-based, with concept representations adaptively formed to  
269 reflect the task at hand<sup>8</sup>. Recent rodent and human work support this proposal with  
270 findings that vmPFC representations are goal-specific in nature, at least at the end of  
271 learning. Specifically, neural ensembles in the rodent homologue of vmPFC have been  
272 demonstrated to represent higher order goal states that relate stimuli to behaviorally-  
273 relevant value<sup>10,25,26</sup>. Similarly, one human neuroimaging study recently localized  
274 representations of a complex task space relating 16 different goal states to vmPFC  
275 activation patterns<sup>7</sup>. Importantly, these vmPFC representations of goal states predicted  
276 participants' behavioral performance, supporting the notion that vmPFC organizes  
277 knowledge based on goals to promote flexible behaviors.

278

279 Our findings provide important evidence for the role of vmPFC during the *formation* of  
280 conceptual maps of experience. Although theoretical perspectives highlight the  
281 importance of vmPFC in cognitive map formation<sup>2,18</sup>, empirical work has failed to directly  
282 examine the computations of vmPFC contributions during encoding. Instead, evidence  
283 is limited to representations that are established after long periods of training<sup>7,12</sup>.  
284 Relatedly, most current models of vmPFC function in memory focus on its role in  
285 biasing reactivation of relevant prior experiences via the hippocampus<sup>27</sup>. Few studies  
286 target the role of vmPFC during encoding. We know that vmPFC interaction with  
287 memory centers during encoding<sup>4,22,28,29</sup>, but we do not know how vmPFC knowledge  
288 representations emerge and support behavior. Our findings provide novel evidence for  
289 vmPFC potentially playing a role in encoding processes that build goal-specific mental  
290 models. The current neural findings are ambiguous as to whether vmPFC is directly  
291 implicated in forming conceptual representations or simply reflects representations  
292 learned elsewhere. For example, rodent models have implicated similar dimensionality-  
293 reduction processes to the basal ganglia with outputs influencing frontal coding<sup>30</sup>.  
294 However, by linking vmPFC coding to the learning mechanisms defined in SUSTAIN,  
295 our results suggest that vmPFC may influence encoding through dimensionality  
296 reduction wherein selective attention highlights goal-specific information and discards  
297 irrelevant dimensions. That vmPFC was the only region identified in our analysis  
298 provides more support for such a direct influence at encoding: inputs to vmPFC are  
299 weighted to select goal-related information and discard irrelevant features in order to  
300 efficiently map input to a goal-directed action. This efficient mapping may then be fed  
301 back to memory centers (i.e., hippocampus) to impact neural coding of learning  
302 experiences<sup>28</sup>. This theorized role for vmPFC coding during learning offers a strong  
303 hypothesis for future work investigating flexible goal-oriented behavior.

304  
305 Our hypothesized view of vmPFC function is based on SUSTAIN's formalism of highly  
306 interactive mechanisms of selective attention and learning<sup>8</sup>, functions theoretically  
307 mapped onto interactions between PFC and the hippocampus<sup>9,28,31</sup>. Support for this  
308 view is found in recent patient work that has demonstrated a causal link between  
309 attentional processes and vmPFC function in decision making<sup>32-34</sup>. These studies have  
310 shown that lesions to vmPFC disrupt attentional guidance based on prior experience  
311 with cue-reward associations<sup>34</sup>, learning the value of task-diagnostic features during  
312 probabilistic learning<sup>32</sup>, and value comparison during reinforcement learning<sup>33</sup>. These  
313 findings have been recently extended to healthy humans in a neuroimaging study which  
314 demonstrated that value signals in vmPFC are dynamically biased by attention during  
315 reinforcement learning<sup>19</sup>. Relatedly, recent rodent work demonstrates the bidirectional  
316 flow of information between the rodent homologue of vmPFC and hippocampus during  
317 context-guided memory encoding and retrieval<sup>35-37</sup>. Coupled with the recent  
318 demonstration of hippocampal-vmPFC functional coupling during concept learning<sup>28</sup>, the

319 current findings align well with the view that vmPFC is critical for evaluating and  
320 representing information in learning and decision making.

321  
322 One limitation of the current work is that we evaluated our neural compression findings  
323 in light of only one computational model. Although SUSTAIN is a well-established model  
324 that explains many learning behaviors<sup>8</sup>, is theoretically motivated by neural mechanisms  
325 of learning<sup>9,31</sup>, and has been successfully linked to neural measures of concept  
326 learning<sup>28,38,39</sup>, an alternative model with similar predictions for attentional tuning over  
327 the course of learning may also account for our neural compression findings. That we  
328 find a significant relationship between SUSTAIN's attentional compression and vmPFC  
329 neural compression suggests the bar for such model comparison is quite high.  
330 However, future studies could leverage our data-driven measure of neural  
331 dimensionality reduction as a target index of learning for adjudicating between formal  
332 cognitive models.

333  
334 The neural compression method proposed here offers a unique approach for evaluating  
335 the informational value of neural representations. One limitation to this approach,  
336 however, is the need for stable trial- or condition-level GLM estimates of BOLD signal.  
337 Such univariate estimates in brain regions with lower signal to noise ratios may be noisy  
338 which would bias PCA towards less dimensionality reduction. Thus, as with any method  
339 that relies on GLM estimation, differences in neural compression are most interpretable  
340 when made across levels of within-participant conditions within the same brain region,  
341 as is the approach in the current study.

342  
343 In summary, we show that learning can be viewed as a process of goal-directed  
344 dimensionality reduction and that such a mechanism is apparent in vmPFC neural  
345 representations throughout learning. Thus, we suggest that vmPFC may play a critical  
346 role not only in representing conceptual content, but in the process of *learning* concepts.  
347 Notably, dimensionality reduction through selective attention offers a reconciling  
348 account of many processes associated with vmPFC including schema representation<sup>40</sup>,  
349 latent casual models<sup>7</sup>, grid-like conceptual maps<sup>12</sup>, and value coding<sup>41,42</sup>.

350

## 351 **Methods**

352

### 353 *Participants*

354

355 Twenty-three volunteers (11 females, mean age 22.3 years old, ranging from 18 to 31  
356 years) participated in the experiment. All subjects were right handed, had normal or  
357 corrected-to-normal vision, and were compensated \$75 for participating.

358

359 The methods used in the current study are novel; however, related category learning  
360 experiments have been employed in several previous studies that focus on analyses of  
361 fMRI activation patterns<sup>43,44</sup>. Given the sample sizes in these studies ( $N = 20, 22, 22$ ),  
362 as well as our previous experience with functional imaging of the whole brain, we set a  
363 target minimum samples size of 20 participants.

364

### 365 *Stimuli*

366

367 Eight color images of insects were used in the experiment (Figure 1A). The insect  
368 images consisted of one body with different combinations of three features: legs, mouth,  
369 and antennae. There were two versions of each feature (thick or thin legs, shovel or  
370 pincer mandible, and thick or thin legs). The eight insect images included all  
371 combination of the three features. The stimuli were sized to 300 x 300 pixels.

372

### 373 *Procedures for the learning problems*

374

375 After an initial screening and consent in accordance with the University of Texas  
376 Institutional Review Board, participants were instructed on the classification learning  
377 problems. Participants then performed the problems in the MRI scanner by viewing  
378 visual stimuli back-projected onto a screen through a mirror attached onto the head coil.  
379 Foam pads were used to minimize head motion. Stimulus presentation and timing was  
380 performed using custom scripts written in Matlab (Mathworks) and Psychtoolbox  
381 ([www.psychtoolbox.org](http://www.psychtoolbox.org)) on an Apple Mac Pro computer running OS X 10.7.

382

383 Participants were instructed to learn to classify the insects based on the combination of  
384 the insects' features using the feedback displayed on each trial. As part of the initial  
385 instructions, participants were made aware of the three features and the two different  
386 values of each feature. Before beginning each classification problem, additional  
387 instructions that described the cover story for the current problem and which buttons to  
388 press for the two insect classes were presented to the participants. One example of this  
389 instruction text is as follows: "Each insect prefers either Warm or Cold temperatures.  
390 The temperature that each insect prefers depends on one or more of its features. On  
391 each trial, you will be shown an insect and you will make a response as to that insect's  
392 preferred temperature. Press the '1' button under your index finger for Warm  
393 temperatures or the '2' button under your middle finger for Cold temperatures." The  
394 other two cover stories involved classifying insects into those that live in the Eastern vs.  
395 Western hemisphere and those that live in an Urban vs. Rural environment. The cover  
396 stories were randomly paired with the three learning problems for each participant. After  
397 the instruction screen, the four fMRI scanning runs (described below) for that problem  
398 commenced, with no further problem instructions. After the four scanning runs for a

399 problem finished, the next problem began with the corresponding cover story  
400 description. Importantly, the rules that defined the classification problems were not  
401 included in any of the instructions; rather, participants had to learn these rules through  
402 trial and error.

403  
404 The three problems the participants learned were structured such that perfect  
405 performance required attending to a distinct set of feature attributes (Figure 1A). For the  
406 low complexity problem, class associations were defined by a rule depending on the  
407 value of one feature attribute. For the medium complexity problem, class associations  
408 were defined by an XOR logical rule that depended on the value of the two feature  
409 attributes that were not relevant in the low complexity problem. For the high complexity  
410 problem, class associations were defined such that all feature attributes had to be  
411 attended to respond correctly. As such, different features were relevant for the three  
412 problems and successful learning required a shift in attending to and representing those  
413 feature attributes most relevant for the current problem. Critically, by varying the number  
414 of diagnostic feature attributes across the three problems, the representational space  
415 for each problem had a distinct informational complexity.

416

stimulus	feature attribute			problem complexity		
	1	2	3	low	medium	high
1	0	0	0	A	A	B
2	0	0	1	A	B	A
3	0	1	0	A	B	A
4	0	1	1	A	A	B
5	1	0	0	B	A	A
6	1	0	1	B	B	B
7	1	1	0	B	B	B
8	1	1	1	B	A	A

417

418 **Table 1:** Stimulus features and class associations for the three learning problems. Each of the eight  
419 stimuli are represented by the binary values of the three feature attributes. The stimuli are assigned to  
420 different classes (A or B) across the low, medium, and high complexity learning problems according to  
421 rules that depend on one, two, or three of the feature attributes, respectively.

422

423 The binary values of the eight insect stimuli along with the class association for the  
424 three learning problems are depicted in Table 1. The stimulus features were randomly  
425 mapped onto the attributes for each participant. These feature-to-attribute mappings  
426 were fixed across the different classification learning problems within a participant. After  
427 the high complexity problem, participants learned the low and medium problems in  
428 sequential order. The learning order of the low and medium problems was

429 counterbalanced across participants. This problem order was used for purposes  
430 described in a prior analysis of this data<sup>28</sup>.

431  
432 The classification problems consisted of learning trials (Figure 1A) during which an  
433 insect image was presented for 3.5s. During stimulus presentation, participants were  
434 instructed to respond to the insect's class by pressing one of two buttons on an fMRI-  
435 compatible button box. Insect images subtended  $7.3^\circ \times 7.3^\circ$  of visual space. The  
436 stimulus presentation period was followed by a 0.5-4.5s fixation. A feedback screen  
437 consisting of the insect image, text of whether the response was correct or incorrect,  
438 and the correct class was shown for 2s followed by a 4-8s fixation. The timing of the  
439 stimulus and feedback phases of the learning trials was jittered to optimize general  
440 linear modeling estimation of the fMRI data. Within one functional run, each of the eight  
441 insect images was presented in four learning trials. The order of the learning trials was  
442 pseudo randomized in blocks of 16 trials such that the eight stimuli were each  
443 presented twice. One functional run was 388s in duration. Each of the learning  
444 problems included four functional runs for a total of 16 repetitions for each insect  
445 stimulus. The entire experiment lasted approximately 65 minutes.

#### 446 447 *Behavioral analysis*

448  
449 Participant-specific learning curves were extracted for each problem by calculating the  
450 average accuracy across blocks of 16 learning trials. These learning curves were used  
451 for the computational learning model analysis. Furthermore, a mixed effect logistic  
452 regression analysis was performed on the behavioral data. Specifically, fixed effects of  
453 stimulus repetition, problem complexity, and learning order along with random intercepts  
454 were estimated in predicting trial-by-trial accuracy across all participants. Accuracy  
455 improved across stimulus repetitions ( $\chi^2 = 769.9$ ,  $p < 0.0001$ ), differed between problem  
456 complexity overall ( $\chi^2 = 970.1$ ,  $p < 0.0001$ ), and changed differently across repetitions  
457 for the problems ( $\chi^2 = 68.9$ ,  $p < 0.0001$ ), but did not differ between learning orders ( $\chi^2 =$   
458  $2.087$ ,  $p < 0.149$ ).

#### 459 460 461 *Computational learning model*

462  
463 Participant behavior was modeled with an established mathematical learning model,  
464 SUSTAIN<sup>8</sup>. SUSTAIN is a network-based learning model that classifies incoming stimuli  
465 by comparing them to memory-based knowledge representations of previously  
466 experienced stimuli. Sensory stimuli are encoded by SUSTAIN into perceptual  
467 representations based on the value of the stimulus features. The values of these  
468 features are biased according to attention weights operationalized as receptive fields on

469 each feature attribute. During learning, these attention weight receptive fields, which  
470 change as a function of the latent model variable  $\lambda_i$ , are tuned to give more weight to  
471 diagnostic features. SUSTAIN represents knowledge as clusters of stimulus features  
472 and class associations that are built and tuned over the course of learning. New clusters  
473 are recruited, and existing clusters updated according to the current learning goals. A  
474 full mathematical formulization of SUSTAIN is provided in its introductory publication<sup>8</sup>.

475  
476 To characterize the attention weights participants formed during learning, we fit  
477 SUSTAIN to each participant's trial-by-trial learning behavior. First, SUSTAIN was  
478 initialized with no clusters and equivalent attention weights across the stimulus feature  
479 attributes. Then, stimuli were presented to SUSTAIN in the same order as a  
480 participant's experience, and model parameters were optimized to predict each  
481 participant's trial-by-trial responses in the three learning problems through a maximum  
482 likelihood optimization method<sup>45</sup>. Specifically, model likelihood was calculated based on  
483 the probability of the model making the same response as the participant in each trial  
484 and this likelihood was maximized through the differential evolution optimization  
485 algorithm provided in the *scipy* python library. In the optimization procedure, the model  
486 state at the end of the first learning problem was used as the initial state for the second  
487 learning problem. In doing so, parameters were optimized to account for learning with  
488 the assumption that attention weights, and knowledge clusters learned from the first  
489 problem carried over to influence learning in the second problem. Similarly, model state  
490 from the second problem carried over and influenced early learning in the third problem.  
491 Thus, problem order effects are considered a natural consequence of our model fitting  
492 approach. The optimized parameters were then used to extract measures of feature  
493 attribute attention weights throughout learning in the three problems. Specifically, for  
494 each participant, the model parameters were fixed to the optimized values and the  
495 model was presented with the trial order experienced by the participant. On each trial,  
496 the values of the feature attribute attention weights,  $\lambda_i$ , were extracted for each  
497 participant. This was repeated for each of the three learning problems. The average  
498 value and 95% confidence intervals of SUSTAIN's five free parameters were:  $\gamma = 8.96 \pm$   
499  $0.82$ ,  $\beta = 1.51 \pm 0.34$ ,  $\eta = 0.08 \pm 0.03$ ,  $d = 17.04 \pm 2.05$ ,  $\tau_h = 0.11 \pm 0.04$ .

#### 500 501 *MRI data acquisition*

502  
503 Whole-brain imaging data were acquired on a 3.0T Siemens Skyra system at the  
504 University of Texas at Austin Imaging Research Center. A high-resolution T1-weighted  
505 MPRAGE structural volume (TR = 1.9s, TE = 2.43ms, flip angle = 9°, FOV = 256mm,  
506 matrix = 256x256, voxel dimensions = 1mm isotropic) was acquired for coregistration  
507 and parcellation. Two oblique coronal T2-weighted structural images were acquired  
508 perpendicular to the main axis of the hippocampus (TR = 13,150ms, TE = 82ms, matrix

509 = 384x384, 0.4x0.4mm in-plane resolution, 1.5mm thru-plane resolution, 60 slices, no  
510 gap). High-resolution functional images were acquired using a T2\*-weighted multiband  
511 accelerated EPI pulse sequence (TR = 2s, TE = 31ms, flip angle = 73°, FOV = 220mm,  
512 matrix = 128x128, slice thickness = 1.7mm, number of slices = 72, multiband factor = 3)  
513 allowing for whole brain coverage with 1.7mm isotropic voxels.

514

#### 515 *MRI data preprocessing and statistical analysis*

516

517 MRI data were preprocessed and analyzed using FSL 5.0.9<sup>46</sup> and custom Python  
518 routines. Functional images were realigned to the first volume of the seventh functional  
519 run to correct for motion, spatially smoothed using a 3mm full-width-half-maximum  
520 Gaussian kernel, high-pass filtered (128s), and detrended to remove linear trends within  
521 each run. Functional images were registered to the MPRAGE structural volume using  
522 Advanced Normalization Tools, version 1.9<sup>47</sup>.

523

#### 524 *Neural compression analysis*

525

526 The goal of the neural compression analysis was to assess the informational complexity  
527 of the neural representations formed during the different learning problems. To index  
528 representational complexity, we measured the extent that neural activation patterns  
529 could be compressed into a smaller dimensional space according to principal  
530 component analyses (PCA). The compression analyses were implemented using  
531 PyMVPA<sup>48</sup> and custom Python routines and were conducted on preprocessed and  
532 spatially smoothed functional data. First, whole brain activation patterns for each  
533 repetition of each stimulus within each run were estimated using an event-specific  
534 univariate general linear model (GLM) approach<sup>16</sup>. This approach allowed us to model  
535 estimates of neural patterns for the eight insect stimuli across the trials in each learning  
536 problem. For each classification problem run, a GLM with separate regressors for  
537 stimulus presentation on each trial, modeled as 3.5s boxcar convolved with a canonical  
538 hemodynamic response function (HRF), was conducted to extract voxel-wise parameter  
539 estimates for each trial. Additionally, trial-specific regressors for the feedback period of  
540 the learning trials (2s boxcar) and responses (impulse function at the time of response),  
541 as well as six motion parameters were included in the GLM. This modeling strategy  
542 targeted the neural representations specific to viewing the stimuli separate from  
543 processes associated with feedback events and trial outcomes for the participants'  
544 responses. This procedure resulted in, for each participant, whole brain activation  
545 patterns for each trial in the three learning problems.

546

547 We assessed the representational complexity of the neural measures of stimulus  
548 representation during learning with a searchlight method<sup>17</sup>. Using a searchlight sphere



549 with a radius of 4 voxels (voxels per sphere: 242 mean, 257 mode, 76 minimum, 257  
550 maximum), we extracted a vector of activation values across all voxels within a  
551 searchlight sphere for all 32 trials within a problem run. These activation vectors were  
552 then submitted to PCA to assess the degree of correlation in voxel activation across the  
553 different trials. PCA was performed using the singular value decomposition method as  
554 implemented in the *decomposition.PCA* function of the *scikit-learn* (version 0.17.1)  
555 Python library. To characterize the amount of dimensional reduction possible in the  
556 neural representation, we calculated the number of principal components that were  
557 necessary to explain 90% of the variance ( $k$ ) in the activation vectors. We scaled this  
558 number into a compression score that ranged from 0 to 1,

559

560

$$compression = 1 - \frac{k}{n},$$

561

562 where  $n$  is equal to 32, the total number of activation patterns submitted to PCA. By  
563 definition, 32 PCs will account for 100% of the variance, but no compression. With this  
564 definition of neural compression, larger compression scores indicated fewer principal  
565 components were needed to explain the variance across trials in the neural data (i.e.,  
566 neural representations with lower dimensional complexity). In contrast, smaller  
567 compression scores indicated more principal components were required to explain the  
568 variance (i.e., neural representations with higher dimensional complexity). This neural  
569 compression searchlight was performed across the whole brain separately for each  
570 participant and each run of the three learning problems in native space. One limitation  
571 that is important to note is that this PCA approach to indexing neural compression does  
572 depend on the success of the single-trial GLM parameter estimates. Brain regions with  
573 lower signal or higher noise may lead to noisy single-trial parameter estimates that  
574 would inflate the PCA estimation of dimensionality (i.e., a higher number of PCs). For  
575 this reason, differences in neural compression are most interpretable when observed for  
576 within-subject factors evaluated within the same brain region, as has been done in the  
577 current work.

578

579 Group-level analyses were performed on the neural compression maps calculated with  
580 the searchlight procedure. Each participant's compression maps were normalized to  
581 MNI space using ANTs<sup>47</sup> and combined into a group dataset. To identify brain regions  
582 that demonstrated neural compression that was consistent with the representational  
583 complexity of the learning problems, we performed a voxel-wise linear mixed effects  
584 regression analysis using the *statsmodels* Python library (version 0.8). The mixed  
585 effects model included factors of problem complexity and learning block as fixed effects  
586 as well as participants as a random effect to predict neural compression. The interaction  
587 of problem complexity and learning block was the central effect of interest. We also  
588 included each participant's accuracy for the three problems within each learning block

589 as a covariate. This regression model was evaluated at each voxel. A statistical map  
590 was constructed by saving the  $t$ -statistic of the interaction between complexity and  
591 learning block. The resulting statistical map was voxel-wise corrected at  $p = 0.001$  and  
592 cluster corrected at  $p = 0.05$  which corresponded to a cluster extent threshold of greater  
593 than 259 voxels. The cluster extent threshold was determined with AFNI<sup>49</sup> 3dClustSim  
594 (version 16.3.12) using the *acf* option, second-nearest neighbor clustering, and 2-sided  
595 thresholding. The 3dClustSim software used was downloaded and compiled on  
596 November 21, 2016 and included fixes for the recently discovered errors of improperly  
597 accounting for edge effects in simulations of small regions and spatial autocorrelation in  
598 smoothness estimates<sup>50</sup>. Additional statistical maps of the main effects of problem  
599 complexity, learning block, and accuracy were also interrogated. No significant clusters  
600 were found for accuracy; see Supplementary Figure 1 and Supplementary Table 1 for  
601 the results for problem complexity and learning block.

602  
603 We assessed the nature of the interaction in the vmPFC cluster by extracting each  
604 participant's average neural compression score within the cluster for each problem  
605 across the four learning runs. This average compression is plotted in the middle panel  
606 of Figure 2B. A Bayesian linear mixed effects model testing the same regression model  
607 as described above was performed on the neural compression scores from the peak  
608 voxel within the vmPFC cluster using the *rstanarm* (version 2.18.2) R library. Relative to  
609 the standard frequentist approach to linear regression, a Bayesian linear mixed effects  
610 approach estimates a full probability model that incorporates uncertainty estimates  
611 about the outcome and predictor variables within a hierarchical framework that explicitly  
612 models participant and group level effects<sup>51</sup>. Through Monte Carlo Markov Chain  
613 (MCMC) procedures, a regression model can be estimated that provides credible  
614 probability estimates for predictor variables without the constraints of normality that limit  
615 frequentist linear regression techniques. The regression models conducted here were  
616 based on default arguments for the *stan\_glm* function: weakly-informed priors with  
617 regularization to prevent over-fitting and four MCMC chains of 2000 samples, the first  
618 half of which are discarded as "warm-up" samples. This results in 4000 total samples  
619 from the posterior distribution of the model. The posterior samples from each factor in  
620 the model can be used to assess model convergence, estimate the average factor  
621 coefficient, define a 95% high-density interval around each factor estimate (i.e., the  
622 Bayesian alternative of confidence intervals), and a  $P$  value representing the proportion  
623 of samples from each factor's posterior distribution that counter the sign of the mean  
624 estimate (i.e., this can be interpreted as a measure of significance similar to frequentist  
625  $p$  values). Model convergence is assessed with the Rhat statistic which estimates the  
626 consistency between independent MCMC chains—values greater than 1.1 suggest the  
627 MCMC sampling did not converge. In all reported Bayesian linear mixed effects model  
628 here, the Rhat values were less than 1.1 suggesting model convergence.

629  
630 It is important to note that separately analyzing neural compression from the peak voxel  
631 within the vmPFC cluster does not represent a set of independent findings. It does,  
632 however, provide a window into the nature of the factors underlying this cluster and the  
633 contribution (or lack thereof) of learning problem order and accuracy. Results from the  
634 Bayesian mixed effects model are summarized in Table 2. The posterior distributions for  
635 learning block, problem complexity, and their interaction are plotted in the right panel of  
636 Figure 2B.

637

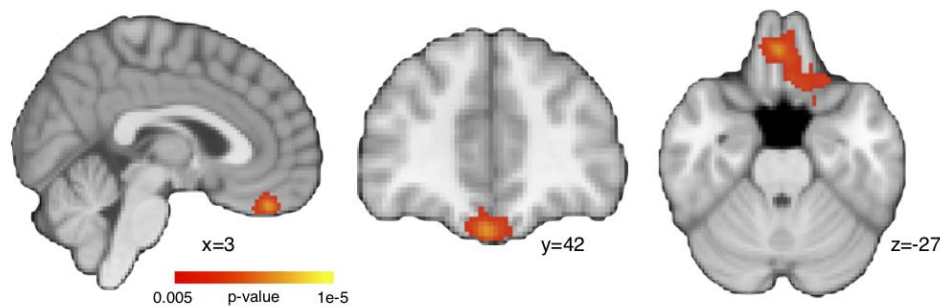
	estimate	95% HDI	<i>P</i>
Intercept	0.586	0.473, 0.699	<0.001
block	0.028	0.014, 0.041	<0.001
complexity	0.025	0.005, 0.044	0.017
accuracy	0.046	0.001, 0.090	0.044
order	0.069	-0.076, 0.209	0.326
block:complexity	-0.013	-0.019, -0.006	<0.001

638  
639 **Table 2:** Results of the Bayesian linear mixed effects regression model predicting neural compression  
640 within the peak voxel of the vmPFC region depicted in Figure 2B. The mean estimated values, 95% high  
641 density interval (HDI), and *P* values are reported for each fixed effect.

642  
643 One potential issue with our findings is that our learning problems vary in their  
644 complexity, and task difficulty can change univariate neural activation. As such, the  
645 concern exists that rather than changes in neural compression, our findings are driven  
646 by simple changes in overall activation levels across problem complexity. To address  
647 this concern, we examined the mean activation in the vmPFC cluster across learning  
648 blocks and problem complexity with a mixed effect linear regression. We found no  
649 differences in average activation across learning block ( $\beta_{mean} = -20.3$ , 95% HDI = [-  
650 134.8, 89.2], *P* = 0.72), problem complexity ( $\beta_{mean} = 25.1$ , 95% HDI = [-115.2, 163.8], *P*  
651 = 0.73), nor an interaction of these factors ( $\beta_{mean} = -1.309$ , 95% HDI = [-82.7, 126.7], *P*  
652 = 0.93). Thus, in our paradigm, problem complexity did not lead to differences in overall  
653 neural activation that changed across learning.

654  
655 One additional concern with our experimental procedure is that the task with the highest  
656 complexity was always learned first. This particular problem order was important for  
657 another purpose in a previous analysis of the data<sup>28</sup>. It is a valid concern that the  
658 differences in neural compression are driven simply by this first problem being the most  
659 difficult and least practiced. Such conditions could potentially be reflected in greater  
660 noise in the neural representations for this higher complexity task which might lead to  
661 lower neural compression, but for reasons not due to problem complexity. To address  
662 this concern, we performed the whole brain voxel-wise linear mixed effects regression

663 analysis comparing complexity across learning blocks while controlling for problem  
664 order and accuracy, but only for the data associated with the low and medium  
665 complexity problems. By excluding the high complexity data, we can target the effect of  
666 complexity without the confound of learning order. If the interaction remains between  
667 problem complexity and learning block in predicting neural compression in the vmPFC  
668 region, it stands to reason that learning order is not a significant driver in the current  
669 findings. Indeed, this follow up analysis revealed a very similar cluster in vmPFC (Figure  
670 5) and no other regions survived cluster correction.  
671



672  
673  
674 **Figure 5:** A whole brain voxel-wise linear mixed effects regression restricted to low and medium  
675 complexity problems (N=23). Similar to the main results in Figure 2B, a cluster in vmPFC showed a  
676 significant interaction between learning block and problem complexity.  
677

678 We also performed a Bayesian linear mixed effects analysis on the peak voxel from the  
679 vmPFC cluster in Figure 2B, but only for the data associated with the low and medium  
680 complexity problems. Again, we found similar results to the full dataset with a significant  
681 interaction between learning block and complexity ( $\beta_{mean} = -0.021$ , 95% HDI = [-0.034, -  
682 0.008],  $P = 0.0015$ ). Both of these analyses suggest that neural compression changes  
683 across learning blocks and increases more for the low versus the medium complexity  
684 problem when restricted to a subset of the data with counterbalanced learning order.  
685

### 686 *Category specific coding in neural compression*

687  
688 The data driven nature of the PCA approach for neural compression can reveal the  
689 degree that neural patterns can be compressed, but not necessarily why this  
690 compression is possible. In the current study, the critical hypothesis is that learning to  
691 attend to problem-specific information will impact neural representations in a manner  
692 consistent with the representational complexity of the problem. Demonstrating that  
693 neural compression increases with learning and does so according to problem  
694 complexity provides compelling evidence in support of our hypothesis, but support that  
695 is nonetheless indirect.  
696

697 To directly assess the contribution of category coding present in the compression  
698 findings, we analyzed how trials for a given problem within a learning block load onto  
699 the identified PCs. If category information is driving neural compression, trials with  
700 stimuli from the same category should load similarly on the PCs and trials from different  
701 categories should load differently on the PCs. In other words the distribution of trial  
702 loadings onto the PCs should discriminate between the two categories. We indexed the  
703 degree of category discrimination in PCA loadings by calculating the absolute difference  
704 between the average loadings within each category for the set of PCs identified in the  
705 neural compression analysis. These category loading differences were weighted by the  
706 explained variance of each PC and summed to create a measure of PC category  
707 discrimination (Figure 4A). Higher values of category discrimination suggests that  
708 category coding is driving neural compression. On the other hand, if trials from both  
709 categories load similarly on the PCs, category discrimination would be equal to 0. A  
710 Bayesian mixed effects linear regression was conducted that evaluated the relationship  
711 between PC category discrimination and factors of learning block, problem complexity,  
712 learning order, and accuracy. We found that category discrimination was present in the  
713 PCA loadings (Figure 4B and Table 3): not only did category discrimination increase  
714 with learning, it did so most for low followed by medium and high complexity problems.  
715

	estimate	95% HDI	<i>P</i>
Intercept	0.0089	-0.0056, 0.0234	0.111
block	0.0080	0.0043, 0.0115	<0.001
complexity	0.0041	-0.001, 0.0089	0.051
accuracy	-0.0029	-0.0132, 0.0079	0.302
order	0.0049	-0.0002, 0.010	0.029
block:complexity	-0.0030	-0.0047, -0.0013	<0.001

716  
717 **Table 3:** Results of the Bayesian linear mixed effects regression model of PC category discrimination  
718 across learning blocks and problem complexity. The mean estimated values, 95% high density interval  
719 (HDI), and *P* values are reported for each fixed effect.

720  
721 *Relating neural compression to behavioral signatures of selective attention*

722  
723 To evaluate the relationship between neural compression and model-based estimates  
724 of attention weighting, we first extracted individual participant-based measures of each.  
725 The participant-specific average neural compression within the mPFC cluster was  
726 extracted for each learning problem. We used the SUSTAIN estimates of stimulus  
727 dimension attention weights,  $\lambda$ , to calculate a signature of selective attention.  
728 Throughout learning on trial-by-trial basis, SUSTAIN tunes attention weights based on  
729 the model parameters, the trial sequence, and the outcome of each trial. For each  
730 participant, we extracted the trial-by-trial derived attention weights in each learning

731 problem based on the participant's best fitting parameters. These attention weights for  
732 the three stimulus dimensions in each problem were transformed to sum to 1, thus  
733 creating a probability distribution representing the likelihood of attention to the three  
734 features. For example, given the attention weights [0.1, 0.1, 0.8], there is a probability of  
735 0.8 that attention will be directed to the third stimulus dimension on any one trial. We  
736 then calculated entropy<sup>38</sup> across the attention weights for each problem separately:

737

$$738 \quad \text{entropy} = -\sum_{i=1}^3 a_i \log_2 a_i,$$

739

740 such that  $a_i$  is the attention weight for stimulus dimension  $i$ . This entropy measure  
741 indexes the dispersion of attention across the stimulus dimensions. If attention is  
742 unselective and all three stimulus dimensions are equally weighted, entropy is high. On  
743 the other hand, if attention is selective with the majority of weight on a single dimension,  
744 entropy is low. To better align attention entropy with our measure of neural  
745 compression, we transformed entropy into an index of attention compression:

746

$$747 \quad \text{attention compression} = 1 - \text{entropy} / \log_2(1/3).$$

748

749 Attention compression first scales entropy according to the maximum amount of entropy  
750 given three stimulus dimensions and then subtracts this ratio from 1. The result is a  
751 measure that ranges from 0 to 1 with low values corresponding to unselective attention  
752 and high values corresponding to more selective attention. As such, the attention  
753 compression index offers a unique signature for optimal attentional strategy across the  
754 three learning problems: the highest attention compression should be seen in the low  
755 complexity problem, an intermediate compression for the medium complexity problem,  
756 and the lowest attention compression for the high complexity problem. As a final step,  
757 for each participant, attention compression was averaged within learning block  
758 separately for each problem. The effect of problem complexity on attention compression  
759 was assessed with a Bayesian linear mixed effects regression that included factors of  
760 learning block, problem complexity, accuracy, learner order, and the interaction of block  
761 and complexity (see Figure 3A and Table 4).

762

	estimate	95% HDI	$P$
Intercept	-0.116	-0.182, -0.049	0.001
block	0.071	0.055, 0.087	<0.001
complexity	0.032	0.009, 0.054	0.006
accuracy	0.103	0.055, 0.151	<0.001
order	-0.025	-0.054, 0.003	0.075
block:complexity	-0.028	-0.035, -0.020	<0.001

763

764 **Table 4:** Results of the Bayesian linear mixed effects regression model predicting attention compression.  
765 The mean estimated values, 95% high density interval (HDI), and  $P$  values are reported for each fixed  
766 effect.

767  
768 We next evaluated the relationship between vmPFC neural compression and attention  
769 weight entropy on an individual participant basis with Bayesian mixed effects linear  
770 regression. The regression model was conducted such that vmPFC neural  
771 compression, learning block, problem complexity, learning order, and accuracy were  
772 predictors of attention weight compression (Figure 3B). This analysis estimates the  
773 degree that vmPFC compression is related to attention compression controlling for all of  
774 the other manipulated factors in the experiment, as well as individual differences in  
775 performance throughout learning. Finding that vmPFC compression is significantly  
776 related to attention compression, in spite of all of these other predictors, suggests that  
777 as attention evolves during learning (according to SUSTAIN's predictions of behavior),  
778 task-specific neural compression is evolving in the same fashion. Indeed, the results  
779 confirm this hypothesis showing a significant correspondence between vmPFC and  
780 attention compression (Table 5). The posterior distribution for the effect of vmPFC  
781 compression on attention compression reveals a robust finding (Figure 3B, right panel).  
782 Importantly, a similar relationship between neural and attention compression was found  
783 when restricting the same regression analyses to only the low and medium complexity  
784 data ( $\beta_{mean} = 0.210$ , 95% HDI = [-0.082, 0.368],  $P = 0.0005$ ).

785

	estimate	95% HDI	$P$
Intercept	-0.214	-0.308, -0.126	<0.001
vmPFC comp.	0.168	0.072, 0.277	<0.001
block	0.067	0.051, 0.083	<0.001
complexity	0.027	0.006, 0.049	0.017
accuracy	0.093	0.045, 0.139	<0.001
Order	-0.031	-0.064, 0.001	0.057
block:complexity	-0.026	-0.033, -0.018	<0.001

786  
787 **Table 5:** Results of the Bayesian linear mixed effects regression model relating attention and vmPFC  
788 compression. The mean estimated values, 95% high density interval (HDI), and  $P$  values are reported for  
789 each fixed effect.

790

#### 791 *Data availability*

792

793 The data collected for this study are available for download: <https://osf.io/5byhb/>.

794

#### 795 *Code availability*

796

797 Data analysis code available upon request.

798

## 799 **Acknowledgments**

800

801 Thanks to Christiane Ahlheim and Margaret Schlichting for manuscript comments. This  
802 work was supported by the Natural Sciences and Engineering Research Council  
803 (Discovery Grant RGPIN-2017-06753 to M.L.M); National Institute of Mental Health  
804 (F32-MH100904 to M.L.M; R01-MH100121 to A.R.P); the Leverhulme Trust (RPG-  
805 2014-075 to B.C.L); the Wellcome Trust (Senior Investigator Award WT106931MA to  
806 B.C.L); and the National Institute of Child Health and Human Development  
807 (1P01HD080679 to B.C.L).

808

## 809 **Author Contributions**

810

811 All authors designed the experiment and wrote the paper. M.L.M. conducted the  
812 research and data analysis.

813

## 814 **Declaration of Interests**

815

816 The authors declare no competing interests.



817 **References**

- 818
- 819 1. Badre, D., Kayser, A. S. & D'Esposito, M. Frontal cortex and the discovery of  
820 abstract action rules. *Neuron* **66**, 315–326 (2010).
  - 821 2. Wilson, R. C., Takahashi, Y. K., Schoenbaum, G. & Niv, Y. Orbitofrontal cortex as  
822 a cognitive map of task space. *Neuron* **81**, 267–278 (2014).
  - 823 3. Mante, V., Sussillo, D., Shenoy, K. V & Newsome, W. T. Context-dependent  
824 computation by recurrent dynamics in prefrontal cortex. *Nature* **503**, 78–84  
825 (2013).
  - 826 4. Zeithamova, D., Dominick, A. L. & Preston, A. R. Hippocampal and Ventral Medial  
827 Prefrontal Activation during Retrieval-Mediated Learning Supports Novel  
828 Inference. *Neuron* **75**, 168–179 (2012).
  - 829 5. Schlichting, M. L., Mumford, J. A. & Preston, A. R. Learning-related  
830 representational changes reveal dissociable integration and separation signatures  
831 in the hippocampus and prefrontal cortex. *Nature Communications* **6**, 8151  
832 (2015).
  - 833 6. Chan, S. C. Y., Niv, Y. & Norman, K. A. A Probability Distribution over Latent  
834 Causes, in the Orbitofrontal Cortex. *The Journal of Neuroscience* **36**, 7817–28  
835 (2016).
  - 836 7. Schuck, N. W. *et al.* Human Orbitofrontal Cortex Represents a Cognitive Map of  
837 State Space Article Human Orbitofrontal Cortex Represents a Cognitive Map of  
838 State Space. *Neuron* **91**, 1402–1412 (2016).
  - 839 8. Love, B. C., Medin, D. & Gureckis, T. M. SUSTAIN: A network model of category  
840 learning. *Psychological Review* **111**, 309–332 (2004).
  - 841 9. Love, B. C. & Gureckis, T. M. Models in search of a brain. *Cognitive, Affective, &*  
842 *Behavioral Neuroscience* **7**, 90–108 (2007).
  - 843 10. Farovik, A. *et al.* Orbitofrontal Cortex Encodes Memories within Value-Based  
844 Schemas and Represents Contexts That Guide Memory Retrieval. *Journal of*  
845 *Neuroscience* **35**, 8333–8344 (2015).
  - 846 11. Zhou, J. *et al.* Rat Orbitofrontal Ensemble Activity Contains Multiplexed but  
847 Dissociable Representations of Value and Task Structure in an Odor Sequence  
848 Task. *Current biology*: *CB* **29**, 897-907.e3 (2019).
  - 849 12. Constantinescu, A. O., O'Reilly, J. X. & Behrens, T. E. J. Organizing conceptual  
850 knowledge in humans with a gridlike code. *Science* **352**, 1464–1468 (2016).
  - 851 13. Leong, Y. C., Radulescu, A., Daniel, R., DeWoskin, V. & Niv, Y. Dynamic  
852 Interaction between Reinforcement Learning and Attention in Multidimensional  
853 Environments. *Neuron* **93**, 451–463 (2017).
  - 854 14. Shepard, R. N., Hovland, C. I. & Jenkins, H. M. Learning and memorization of  
855 classification. *Psychological Monographs* **75**, 517 (1961).
  - 856 15. Eckart, C. & Young, G. The Approximation of One Matrix by Another Low Rank.  
857 *Psychometrika* **1**, 211–218 (1936).
  - 858 16. Mumford, J. A., Turner, B. O., Ashby, F. G. & Poldrack, R. A. Deconvolving BOLD  
859 activation in event-related designs for multivoxel pattern classification analyses.  
860 *NeuroImage* **59**, 2636–2643 (2012).

- 861 17. Kriegeskorte, N., Goebel, R. & Bandettini, P. Information-based functional brain  
862 mapping. *Proceedings of the National Academy of Sciences of the United States*  
863 *of America* **103**, 3863–3868 (2006).
- 864 18. Wikenheiser, A. M. & Schoenbaum, G. Over the river, through the woods:  
865 cognitive maps in the hippocampus and orbitofrontal cortex. *Nat Rev Neurosci* **17**,  
866 513–523 (2016).
- 867 19. Leong, Y. C., Radulescu, A., Daniel, R., DeWoskin, V. & Niv, Y. Dynamic  
868 Interaction between Reinforcement Learning and Attention in Multidimensional  
869 Environments. *Neuron* **93**, 451–463 (2017).
- 870 20. Gilboa, A. & Marlatt, H. Neurobiology of Schemas and Schema-Mediated  
871 Memory. *Trends in Cognitive Sciences* **14**, 417–428 (2017).
- 872 21. Tse, D. *et al.* Schema-Dependent Gene Activation. *Science* **333**, 891–895 (2011).
- 873 22. van Kesteren, M. T. R., Fernández, G., Norris, D. G. & Hermans, E. J. Persistent  
874 schema-dependent hippocampal-neocortical connectivity during memory  
875 encoding and postencoding rest in humans. *Proceedings of the National*  
876 *Academy of Sciences of the United States of America* **107**, 7550–7555 (2010).
- 877 23. Spalding, K. N., Jones, S. H., Duff, M. C., Tranel, D. & Warren, D. E. Investigating  
878 the Neural Correlates of Schemas: Ventromedial Prefrontal Cortex Is Necessary  
879 for Normal Schematic Influence on Memory. *J Neurosci* **35**, 15746–15751 (2015).
- 880 24. Ghosh, V. E., Moscovitch, M., Melo Colella, B. & Gilboa, A. Schema  
881 representation in patients with ventromedial PFC lesions. *The Journal of*  
882 *Neuroscience* **34**, 12057–70 (2014).
- 883 25. Lopatina, N. *et al.* Ensembles in medial and lateral orbitofrontal cortex construct  
884 cognitive maps emphasizing different features of the behavioral landscape.  
885 *Behavioral Neuroscience* **131**, 201–212 (2017).
- 886 26. Zhou, J. *et al.* Rat Orbitofrontal Ensemble Activity Contains Multiplexed but  
887 Dissociable Representations of Value and Task Structure in an Odor Sequence  
888 Task. *Current biology*: *CB* **29**, 897–907.e3 (2019).
- 889 27. Miller, E. K. & Cohen, J. D. An integrative theory of prefrontal cortex function.  
890 *Annual Review of Neuroscience* **24**, 167–202 (2001).
- 891 28. Mack, M. L., Love, B. C. & Preston, A. R. Dynamic updating of hippocampal  
892 object representations reflects new conceptual knowledge. *Proceedings of the*  
893 *National Academy of Sciences* **113**, 13203–13208 (2016).
- 894 29. Schlichting, M. L. & Preston, A. R. Hippocampal–medial prefrontal circuit supports  
895 memory updating during learning and post-encoding rest. *Neurobiology of*  
896 *Learning and Memory* **134**, 91–106 (2016).
- 897 30. Bar-Gad, I., Havazelet-Heimer, G., Goldberg, J. A., Ruppin, E. & Bergman, H.  
898 Reinforcement-driven dimensionality reduction--a model for information  
899 processing in the basal ganglia. *Journal of basic and clinical physiology and*  
900 *pharmacology* **11**, 305–20 (2000).
- 901 31. Mack, M. L., Love, B. C. & Preston, A. R. Building concepts one episode at a  
902 time: The hippocampus and concept formation. *Neuroscience Letters* **680**, 31–38  
903 (2018).
- 904 32. Vaidya, A. R. & Fellows, L. K. Necessary Contributions of Human Frontal Lobe

- 905 Subregions to Reward Learning in a Dynamic, Multidimensional Environment. *The*  
906 *Journal of Neuroscience* **36**, 9843–58 (2016).
- 907 33. Noonan, M. P., Chau, B. K. H., Rushworth, M. F. S. & Fellows, L. K. Contrasting  
908 Effects of Medial and Lateral Orbitofrontal Cortex Lesions on Credit Assignment  
909 and Decision-Making in Humans. **37**, 7023–7035 (2017).
- 910 34. Vaidya, A. R. & Fellows, L. K. Ventromedial Frontal Cortex Is Critical for Guiding  
911 Attention to Reward-Predictive Visual Features in Humans. *Journal of*  
912 *Neuroscience* **35**, (2015).
- 913 35. Place, R., Farovik, A., Brockmann, M. & Eichenbaum, H. Bidirectional prefrontal-  
914 hippocampal interactions support context-guided memory. *Nature Neuroscience*  
915 **19**, (2016).
- 916 36. Wikenheiser, A. M., Marrero-Garcia, Y. & Schoenbaum, G. Suppression of  
917 Ventral Hippocampal Output Impairs Integrated Orbitofrontal Encoding of Task  
918 Structure. *Neuron* **95**, 1197-1207.e3 (2017).
- 919 37. Guise, K. G. & Shapiro, M. L. Medial Prefrontal Cortex Reduces Memory  
920 Interference by Modifying Hippocampal Encoding Article Medial Prefrontal Cortex  
921 Reduces Memory Interference by Modifying Hippocampal Encoding. *Neuron* **94**,  
922 183-192.e8 (2017).
- 923 38. Davis, T., Love, B. C. & Preston, A. R. Striatal and hippocampal entropy and  
924 recognition signals in category learning: Simultaneous processes revealed by  
925 model-based fMRI. *Journal of Experimental Psychology: Learning, Memory, and*  
926 *Cognition* **38**, 821–839 (2012).
- 927 39. Davis, T., Love, B. C. & Preston, A. R. Learning the exception to the rule: Model-  
928 based fMRI reveals specialized representations for surprising category members.  
929 *Cerebral Cortex* **22**, 260–273 (2012).
- 930 40. Van Kesteren, M. T. R., Ruiters, D. J., Fernández, G. & Henson, R. N. How  
931 schema and novelty augment memory formation. *Trends in Neurosciences* **35**,  
932 211–219 (2012).
- 933 41. Clithero, J. A. & Rangel, A. Informatic parcellation of the network involved in the  
934 computation of subjective value. *Social Cognitive and Affective Neuroscience* **9**,  
935 1289–1302 (2013).
- 936 42. Grueschow, M., Polania, R., Hare, T. A. & Ruff, C. C. Automatic versus Choice-  
937 Dependent Value Representations in the Human Brain. *Neuron* **85**, 874–885  
938 (2015).
- 939 43. Mack, M. L., Preston, A. R. & Love, B. C. Decoding the brain's algorithm for  
940 categorization from its neural implementation. *Current Biology* **23**, 2023–2027  
941 (2013).
- 942 44. Davis, T., Xue, G., Love, B. C., Preston, A. R. & Poldrack, R. A. Global neural  
943 pattern similarity as a common basis for categorization and recognition memory.  
944 *The Journal of neuroscience*: the official journal of the Society for Neuroscience  
945 **34**, 7472–84 (2014).
- 946 45. Storn, R. & Price, K. Differential evolution—a simple and efficient heuristic for  
947 global optimization over continuous spaces. *Journal of Global Optimization* **11**,  
948 341–359 (1997).

- 949 46. Jenkinson, M., Beckmann, C. F., Behrens, T. E. J., Woolrich, M. W. & Smith, S.  
950 M. FSL. *NeuroImage* **62**, 782–90 (2012).
- 951 47. Avants, B. B. *et al.* A reproducible evaluation of ANTs similarity metric  
952 performance in brain image registration. *NeuroImage* **54**, 2033–2044 (2011).
- 953 48. Hanke, M. *et al.* PyMVPA: A python toolbox for multivariate pattern analysis of  
954 fMRI data. *Neuroinformatics* **7**, 37–53 (2009).
- 955 49. Cox, R. W. AFNI: software for analysis and visualization of functional magnetic  
956 resonance neuroimages. *Computers and Biomedical Research* **29**, 162–173  
957 (1996).
- 958 50. Eklund, A., Nichols, T. E. & Knutsson, H. Cluster failure: Why fMRI inferences for  
959 spatial extent have inflated false-positive rates. *Proceedings of the National  
960 Academy of Sciences* **113**, 201602413 (2016).
- 961 51. Muth, C., Oravecz, Z. & Gabry, J. User-friendly Bayesian regression modeling: A  
962 tutorial with rstanarm and shinystan. **14**, (2018).
- 963