

Occipitotemporal Representations Reflect Individual Differences in Conceptual
Knowledge

Kurt Braunlich

University College London

Bradley C. Love

University College London, The Allan Turing Institute

Total Words in text: 8656

Author Note

We thank the authors of the original studies for sharing their data. This work was supported by NIH Grant 1P01HD080679, Leverhulme Trust grant RPG-2014-075 and Wellcome Trust Senior Investigator Award WT106931MA to BCL.

Abstract

Through selective attention, decision-makers can learn to ignore behaviorally-irrelevant stimulus dimensions. This can improve learning and increase the perceptual discriminability of relevant stimulus information. Across cognitive models of categorization, this is typically accomplished through the inclusion of attentional parameters, which provide information about the importance assigned to each stimulus dimension by each participant. The effect of these parameters on psychological representation is often described geometrically, such that perceptual differences over relevant psychological dimensions are accentuated (or stretched), and differences over irrelevant dimensions are down-weighted (or compressed). In sensory and association cortex, representations of stimulus features are known to covary with their behavioral relevance. Although this implies that neural representational space might closely resemble that hypothesized by formal categorization theory, to date, attentional effects in the brain have been demonstrated through powerful experimental manipulations (e.g., contrasts between relevant and irrelevant features). This approach sidesteps the role of idiosyncratic conceptual knowledge in guiding attention to useful information sources. To bridge this divide, we used formal categorization models, which were fit to behavioral data, to make inferences about the concepts and strategies used by individual participants during decision-making. We found that when greater attentional weight was devoted to a particular visual feature (e.g., “color”), its value (e.g., “red”) was more accurately decoded from occipitotemporal cortex. We additionally found that this effect was sufficiently sensitive to reflect individual differences in conceptual knowledge, indicating that occipitotemporal stimulus representations are embedded within a space closely resembling that formalized by classic categorization theory.

Keywords: Concepts, Selective Attention, Occipitotemporal Cortex

Occipitotemporal Representations Reflect Individual Differences in Conceptual Knowledge

Introduction

Through selective attention, knowledge of abstract concepts can emphasize relevant stimulus features. For example, while the size of garments is critical when choosing what to purchase, weight may be more important when deciding how to ship them. The attention devoted to individual features is flexibly-modulated according to current goals, transient contextual demands, and reflects evolving conceptual knowledge (Goldstone, 2003; Tversky, 1977). In formal categorization models, a way to account for this flexibility is through inclusion of attentional parameters, which reflect the influence of each dimension on the category decision (e.g., Kruschke, 1992; Love, Medin, & Gureckis, 2004; Nosofsky, 1986). These attentional parameters are often described as “warping” multidimensional psychological space, such that differences along relevant stimulus dimensions are accentuated (or “stretched”) and differences along irrelevant dimensions are down-weighted (or “compressed”; Figure 1). Here, we directly test this classic idea by investigating whether the strength of neural stimulus feature representations are modulated by these attentional parameters. Importantly, we attempt to relate individual differences in conceptual knowledge (as revealed by model fits of attentional parameters) to individual differences in neural representation (as revealed by decoding stimulus features in fMRI data). In doing so, we aim to bridge behavioral and neural levels of analysis at the individual level using cognitive models.

Figures

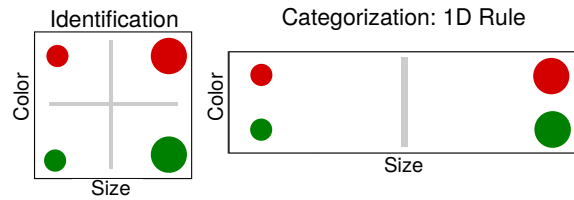


Figure 1. Example: Attention Influences Psychological Space. Left: In an object identification task, both psychological dimensions should receive equivalent attention, as they are equally relevant. **Right:** In a one-dimensional rule-based categorization task, only a single dimension is relevant (in this example, size), and decision-makers could ignore the irrelevant dimension (color). This is often described as “warping” psychological space such that differences along relevant dimensions are accentuated (or “stretched”), and differences along irrelevant dimensions are down-weighted (or “compressed”).

When identifying specific objects, agents must typically consider all stimulus features, and the psychological distance between stimuli closely reflects their perceptual attributes (Shepard, 1957; Townsend & Ashby, 1982). During categorization, however, groups of distinct stimuli must be treated equivalently, and both learning and generalization can be improved by selectively attending to relevant stimulus dimensions (Nosofsky, 1986; Shepard, Hovland, & Jenkins, 1961). Although categorization models differ in how stimuli are represented in memory (e.g., as individual exemplars, as prototypes, or as clusters that flexibly reflect environmental structure; Love et al., 2004; Minda & Smith, 2002; Nosofsky, 1987; Nosofsky & Zaki, 2002; Smith & Minda, 1998; Zaki, Nosofsky, Stanton, & Cohen, 2003), they similarly assume that categorization involves learning to distribute attention across stimulus features so as to optimize behavioral performance. Although they differ in their mathematical details, these models also posit that *endogenous* (i.e., “top-down”) attentional control (Miller & Cohen, 2001; Tsotsos, 2011) can modulate the influence of the *exogenous* (or perceptual) stimulus dimensions on the behavioral choice. The attentional parameters play a key role in allowing the models to capture patterns of human generalization across different goals and different rules. As they also predict human eye-movements during category decision-making (e.g., Rehder & Hoffman, 2005a, 2005b), they are thought to reflect the strategies used by individual decision-makers to integrate

information from the external world.

In the brain, effects of endogenous attention have been observed across the visual cortical hierarchy (Buffalo, Fries, Landman, Liang, & Desimone, 2010; Jehee, Brady, & Tong, 2011; Kamitani & Tong, 2005, 2006; Luck, Chelazzi, Hillyard, & Desimone, 1997; Motter, 1993). A general finding is that when attention is devoted to a specific visual feature, its neural representation is more accurately decoded. For instance, in human fMRI, when multiple visual gratings are concurrently presented, representations of attended orientations in areas V1-V4 are more easily decoded than those that are unattended (Jehee et al., 2011; Kamitani & Tong, 2005). Similarly, when random dot stimuli move in multiple directions, representations of attended motion directions in area MT+ are more easily decoded than those that are unattended (Kamitani & Tong, 2006). Whereas these studies have relied on explicit cues to guide attention to relevant aspects of the stimulus array, in real-world environments, decision-makers must typically rely on knowledge gained through past experience in order to selectively attend to relevant information sources.

Categorization tasks mirror this aspect of real-world environments; decision-makers must rely on learned conceptual knowledge in order to selectively attend to relevant stimulus dimensions. Several studies have investigated whether neural representations of exogenous information sources are modulated by learned conceptual knowledge (e.g., Folstein, Palmeri, & Gauthier, 2013; Li, Ostwald, Giese, & Kourtzi, 2007; Sigala & Logothetis, 2002). Sigala and Logothetis (2002), for instance, trained macaques to categorize abstract images, which varied according to four stimulus dimensions. Neural representations of the two behaviorally-relevant stimulus dimensions (i.e., the dimensions that reliably predicted the correct response) in the inferior temporal lobe were enhanced relative to those of the irrelevant dimensions. Using fMRI with human participants, Li et al. (2007) investigated whether neural representations of stimulus motion and shape were influenced by their relevance to the active categorization rule. Using multivariate pattern analysis (MVPA), they similarly found that representations of these stimulus dimensions reflected their relevance to the active

rule.

Across studies involving explicit attentional cues and categorization studies involving learned conceptual knowledge, a general finding is that occipitotemporal representations of behaviorally-relevant information sources are enhanced relative to those that are irrelevant (this may not hold for integral stimulus dimensions; Garner, 1976). These effects are compelling, as they imply that occipitotemporal representational space may closely resemble that conceptualized by classic cognitive theory (e.g., Kruschke, 1992; Love et al., 2004; Nosofsky, 1986). Specifically, it may expand and contract, along axes defined by perceptually-separable stimulus dimensions (Garner, 1976), in ways that closely reflect the idiosyncratic concepts and strategies used by individual participants during decision-making.



Previous studies have relied on contrastive analyses, in which neural representations of attended stimulus dimensions are compared to those of unattended dimensions. Although statistically-powerful, this approach defines selective attention in terms of the experimental paradigm (but see O'Bryan, Walden, Serra, & Davis, 2018), and therefore sidesteps effects associated with individual differences in conceptual knowledge (e.g., Craig & Lewandowsky, 2012; Little & McDaniel, 2015; McDaniel, Cahill, Robbins, & Wiener, 2014; Raijmakers, Schmittmann, & Visser, 2014). These effects can be substantial, particularly for ill-defined categorization-problems (such as the 5/4 categorization task), which are common in every-day life (Hedge, Powell, & Sumner, 2017; Johansen & Palmeri, 2002). Here, we bridge this divide by combining model-based fMRI (Palmeri, Love, & Turner, 2017; Turner, Forstmann, Love, Palmeri, & Van Maanen, 2017) with multivariate pattern analyses. This allowed us to abstract away from individual differences in neural topography (Haxby et al., 2001; Haynes, 2015; Kriegeskorte & Kievit, 2013), to investigate whether neural stimulus representations reflect individual differences in conceptual knowledge. Specifically, we sought to investigate whether the attentional parameters derived from formal categorization models predict contortions of occipitotemporal representational space during decision-making.

We investigated this hypothesis using two publicly available datasets (osf.io). In the first (Mack, Preston, & Love, 2013), participants categorized abstract stimuli that varied according to four binary dimensions (Figure 2.A), according to a categorization strategy they learned prior to scanning. In the original paper, the authors fit both the Generalized Context Model (GCM; Nosofsky, 1986) and the Multiplicative Prototype Model (Nosofsky, 1987; Nosofsky & Zaki, 2002) to the behavioral data, and used them to compare exemplar and prototype accounts of occipitotemporal representation. Using representational similarity analysis (Kriegeskorte, Mur, & Bandettini, 2008), Mack et al. (2013) additionally identified regions of the brain (lateral occipital cortex, parietal cortex, inferior frontal gyrus, and insular cortex) sensitive to the attentionally-modulated pairwise similarities between stimuli. Although these results (particularly those in lateral occipital cortex) imply that neural representations of the individual stimulus features might be modulated by selective attention, in principle, this could also reflect modulation within an abstract representational space where stimulus features are not individually represented. For instance, while visual cortex reflects sensory input (and is known to represent individual stimulus dimensions), prefrontal cortex can flexibly represent conjunctions of features, abstract rules, and category boundaries in a goal-directed manner. Representations in parietal cortex display intermediate characteristics, as they can reflect both sensory and decisional factors (Brincat, Siegel, Nicolai, & Miller, 2017; Jiang et al., 2007; Li et al., 2007).

In the second dataset (Mack, Love, & Preston, 2016), participants learned, while scanning, to categorize images of insects that varied according to three binary perceptual dimensions (Figure 2.B), according to type I, type II and type VI problems described by Shepard et al. (1961).¹ Importantly, although the same stimuli were included in each task, the degree to which each of the features predicted the correct choice differed between rules. The authors fit the SUSTAIN learning model (Supervised and Unsupervised Stratified Adaptive Incremental Network; Love et al., 2004) to the behavioral data, and used it to investigate hippocampal involvement in the development

¹In their paper, Mack et al. (2016) focus on effects associated with the type I and type II rules.

A) 5/4 Experiment

		
Color:	red	green
Size:	large	small
Shape:	circle	triangle
Position:	left	right

B) SHJ Experiment



		
Mandible:	pincer	shovel
Antennae:	thick	thin
Legs:	thin	thick

Figure 2. Stimuli. **A)** Two of the 16 stimuli used in the “5/4” experiment are illustrated. The stimuli varied according to four binary perceptual dimensions: color, size, shape and position. **B)** Two of the eight stimuli used in the SHJ experiment are illustrated. The stimuli were pictures of insects that varied according to three binary dimensions: mandible shape (highlighted in green), antennae thickness (highlighted in blue), and leg thickness (highlighted in red). For both experiments, the mapping of visual dimension to its role in each category structure (Tables 1 and 2) was randomized for each participant.

of new conceptual knowledge. Using representational similarity analysis, they found that SUSTAIN successfully predicted the pairwise similarities between hippocampal stimulus representations across rule-switches. This suggests that hippocampal representations are updated according to goal-directed attentional selection of stimulus features.

Methods

Description of Datasets

In both experiments, participants categorized stimuli that were characterized by multiple perceptually-separable stimulus dimensions. As the mapping of perceptual attributes to their role in each category structure was randomized for each participant, it is possible to differentiate effects associated with intrinsic perceptual stimulus attributes from effects of behavioral relevance. For example, while color strongly

predicted the correct category choice for some participants, it provided unreliable information for others. In both experiments, participants were not instructed as to which cues were informative, and learned to perform each task through trial-and-error.

We used the GCM for the first dataset (the winning model from Mack et al. 2013), as participants learned how to perform the categorization task prior to scanning. We used SUSTAIN for the second dataset, as it learns on a trial-by-trial basis, and participants learned to perform each task during scanning. SUSTAIN was additionally fit in such a way that the learning of one task carried over to the next. Importantly, although the GCM and SUSTAIN differ in how stimuli are represented in memory (i.e., as exemplars or clusters), they similarly posit that attention “contorts” psychological space, as illustrated in Figure 1. Thus, these studies and models provide a good test of whether attention weights in successful cognitive models are plausible at both behavioral and neural levels of analysis.

The “5/4” Dataset. The first dataset (Mack et al., 2013) was collected while 20 participants (14 Female) categorized abstract stimuli (Figure 2.A), which varied according to four binary stimulus dimensions (size: large vs. small, shape: circle vs. triangle, color: red vs. green, and position: left vs. right). Prior to scanning, they learned to categorize the stimuli according to the “5/4” categorization task (Medin & Schaffer, 1978) through trial-and-error. During this training session, participants were shown only the first nine stimuli shown in Table 1 (i.e., five category “A” members: A1-A5 and four category “B” members: B1-B4), and experienced 20 repetitions of each stimulus. During the anatomical scan, they additionally performed a “refresher” task, involving four additional repetitions on each training item. Each training trial involved a 3.5 second stimulus presentation period in which participants made a button press. Following the button press, a fixation cross was shown for 0.5 seconds, and feedback was then presented for 3.5 seconds. Feedback included information about the correct category, and about whether the response was correct or incorrect. During scanning, participants were required to categorize not only the training items, but also the seven transfer stimuli (i.e., T1-T7). In the scanner, stimuli were presented for 3.5 seconds on

each trial, no feedback was provided, and stimuli were separated by a 6.5 second intertrial interval. Over six runs, each of the 16 stimuli were presented three times. The order of the stimulus presentations were randomized for each participant.

Stimulus	D1	D2	D3	D4
A1	1	0	0	0
A2	1	0	1	0
A3	0	1	0	0
A4	0	0	1	0
A5	0	0	0	1
B1	1	1	0	0
B2	1	0	0	1
B3	0	1	1	1
B4	1	1	1	1
T1	0	1	1	0
T2	1	1	1	0
T3	0	0	0	0
T4	1	1	0	1
T5	0	1	0	1
T6	0	0	1	1
T7	1	0	1	1

Table 1

The “5/4” Category Structure. *Prior to scanning, participants learned, through trial and error, to categorize the first nine stimuli (category “A”: A1-A5; category “B”: B1-B4) illustrated in Figure 2.A. During scanning, they categorized both the training and the transfer (T1-T7) stimuli. Perceptual stimulus dimensions (Figure 2) were pseudo-randomly assigned to category dimensions for each participant.*

Whole-brain images were acquired 3T GE Medical Systems Signa scanner. Structural images were collected using a T2-weighted, flow compensated spin-echo pulse sequence (TR=3s, TE=68ms, 256×256 matrix, 1×1mm in-plane resolution, 33 slices, 3mm slice thickness, gap=0.6mm). An additional T1-weighted 3D SPGR structural image was also collected (256×256×172 matrix, 1×1×1.3mm voxels). Functional images were collected using an echo planar imaging sequence (TR=2s, TE=30.5ms, flip angle=73°, 64×64 matrix, 3.75×3.75 in-plane resolution, bottom-up interleaved sequence, gap=0.6mm).

The SHJ Dataset. In the second dataset (Mack et al., 2016), 23 right-handed participants (11 Female, mean age = 22.3 years) categorized images of insects (Figure

2.B) varying along three binary dimensions (legs: thick vs. thin, antennae: thick vs. thin, and mandible: pincer vs. shovel). We excluded data from two participants who each had corrupted data on one run. This resulted in 21 participants for the final analyses. During scanning, participants learned to categorize the stimuli according to the type I, type II, and type VI problems described by Shepard et al. (1961). In the type I problem, the optimal strategy required attending to a single stimulus dimension (e.g., “legs”) that perfectly predicted the category label, while ignoring the other two dimensions. In the type II problem, the optimal strategy was a logical XOR rule, in which two stimulus features had to be considered together. In the type VI problem, all stimulus features were relevant to the decision, and participants had to learn the mapping between individual stimuli and the category label. To maximally differentiate endogenous and exogenous factors, the irrelevant feature in the type II rule was used as a relevant feature of the type I problem for each participant.

Stimulus	D1	D2	D3	Type1	Type2	Type6
1	0	0	0	A	A	A
2	0	0	1	A	B	B
3	0	1	0	A	B	B
4	0	1	1	A	A	A
5	1	0	0	B	A	B
6	1	0	1	B	B	A
7	1	1	0	B	B	A
8	1	1	1	B	A	B

Table 2

SHJ Category Structures. Participants learned by trial-and-error to perform the type I (a one-dimensional rule-based categorization task), type II (a two-dimensional XOR rule-based categorization task), and type VI (a three-dimensional task requiring memorization of the individual stimuli) problems during scanning. For each participant, perceptual stimulus dimensions (Figure 2) were randomly assigned to these abstract category dimensions.

Each problem was performed across four scanner runs. While all of the participants learned to perform the type VI problem first, the order of the type I and type II problems was then counterbalanced across participants. Each trial consisted of a 3.5 second stimulus presentation period, a jittered 0.5-4.5 second fixation period, and

feedback. Feedback was presented for 2 seconds and consisted of an image of the presented insect, as well as text indicating whether the response was correct or incorrect. Each trial was separated by jittered intertrial interval (4-8 seconds), which consisting of a fixation cross. Each run included four presentations of each of the eight stimuli.

For consistency across datasets, we used the group-derived region of interest (ROI) used in “5/4” dataset (Figure 3.B), and performed a similar analysis. As participants in the SHJ experiment learned to perform the type I, type II and type VI problems during scanning, we mirrored the strategy used by the original authors, and divided the scanning sessions into early (first two runs of each problem) and late learning epochs (last two runs of each problem). We investigated the relationship between occipitotemporal representation and attention only during this late learning phase, in which behavior had largely stabilized.

Whole-brain images were acquired in a 3T Siemens Skyra Scanner. Anatomical images were collected using a T1-weighted MPRAGE sequence (TR=1.9s, TE=2.43ms, 256×256 matrix, 1mm isotropic voxels, flip angle=9°, FOV=256mm). Functional images were acquired using a T2*-weighted multiband (multiband factor=3) accelerated EPI sequence (TR=2s, TE=31ms, flip angle=73°, FOV=220mm, 128×128 matrix, 1.7mm slice thickness, 1.7mm isotropic voxels).

SUSTAIN was initialized with no clusters, and with equivalent weights assigned to each stimulus dimension. Its learning parameters were first fit to the learning performance of each participant using a maximum-likelihood genetic algorithm procedure. The model was fit in such a way that, after learning one problem, the model state was used as the initial state for the subsequent problem. In this way, the model was fit under the assumption that learning of one task would influence later behavior. Once the learning parameters of the model were optimized, they were fixed, and the attentional parameters were extracted from the second two runs of each task (in which learning had largely stabilized). This yielded distinct sets of attentional parameters for each participant and each task. More information about the model can be found in appendix B.

Image Processing

Preprocessing included motion correction, and coregistration of the anatomical images to the mean of the functional images (using Statistical Parametric Mapping (SPM), version 6470). All MVPA analyses were performed in native space without smoothing. For group-level analyses, the statistical maps from each participant were warped to Montreal Neurological Institute (MNI) atlas space using Advanced Normalization Tools (ANTs; Avants, Tustison, & Johnson, 2009), and then smoothed with a 6mm full-width at half maximum Gaussian kernel. The ROI derived from group-level analyses were transformed back into each participants native space for ROI-level analyses. We performed MVPA on the unsmoothed, single-trial, *t*-statistic images (Misaki, Kim, Bandettini, & Kriegeskorte, 2010) derived from the least-squares separate procedure (LSS; Mumford, Turner, Ashby, & Poldrack, 2012). We used SPM to estimate the LSS images for the “5/4” dataset, but used the NiPy python package (<http://nipy.org/nipy/index.html>) for the SHJ dataset, as it tends to run more efficiently, and this study used a multiband sequence with smaller voxel dimensions.

Results

Representations of Individual Visual Features

To identify regions most strongly representing the stimulus features, we performed a cross-validated searchlight analysis (sphere radius = 10mm.; Kriegeskorte, Goebel, & Bandettini, 2006)² in which we decoded each of the four visual features (position, shape, color and size). We performed the analysis in native anatomical space, using a linear support vector classifier (SVC, C=0.1; using the Scikit-learn python package; Pedregosa et al., 2011) in conjunction with a five-fold, leave-one-run out, cross-validated procedure. This involved repeatedly training the model on four of the five runs, and testing whether it could accurately predict the stimulus features associated with the held-out neuroimaging data.

²This involves moving an imaginary sphere throughout the brain; repeatedly investigating how well the voxels within the sphere can decode a variable of interest.

After centering each of the resultant statistical maps at chance (50% for each visual feature), we created a single map for each participant, which reflected the average, above chance, decoding accuracy across features. We then normalized each map to MNI space and, in order to identify regions supporting above chance feature decoding, performed a group-level permutation test. This involved randomly flipping the sign of the statistical maps 10,000 times (using the `randomise` function from the Oxford Centre for Functional MRI of the Brain Software Library (FSL); Winkler, Ridgway, Webster, Smith, & Nichols, 2014). The familywise error rate was controlled using a voxelwise threshold of $p < 0.001$. This identified right middle frontal gyrus (BA9) and left post-central motor cortex, as well as widespread visual and association cortex, extending dorsally from occipital pole to the bilateral superior extrastriate cortex and bilateral intraparietal sulcus (IPS), and ventrally into the bilateral lingual gyrus (Table 3). As this procedure yielded a diffuse pattern of spatial activity, we increased the minimum t -statistic threshold (from 6.24 to 9) to isolate voxels most strongly representing the individual stimulus features. This removed voxels belonging to the bilateral inferior occipital cortex, left lingual gyrus, bilateral intraparietal sulcus, and bilateral precuneus. The resultant ROI is illustrated in Figure 3.B.

Effects Associated with Conceptual Knowledge

“5/4” Dataset. First, we confirmed that each stimulus feature could be decoded significantly above chance from the ROI illustrated in Figure 3.B. Although estimating effect sizes on voxels selected through non-orthogonal criteria is circular, testing significance at the ROI-level has been recommended to confirm that information exists, not only at the level of the searchlight sphere, but also at the level of the ROI (Etzel, Zacks, & Braver, 2013). This analysis also allows us to illustrate the individual feature decoding accuracies for each participant (Figure 3.C). The analyses were performed in the native anatomical space of each participant using the cross-validated SVC analysis described above (but setting the C parameter to 1 instead of 0.1, which

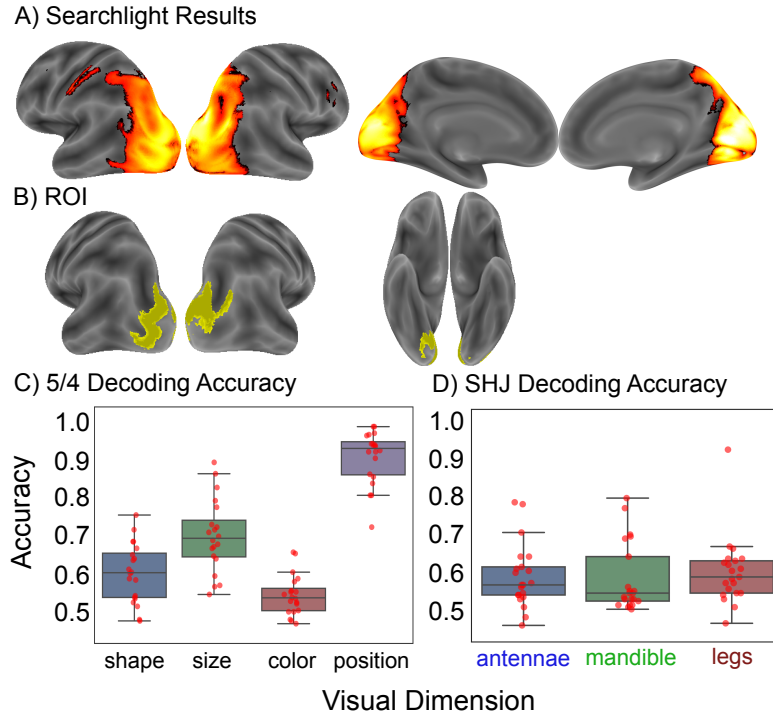


Figure 3. **A)** For the “5/4” dataset, a searchlight analysis indicated that binary perceptual dimensions could be decoded from widespread visual regions (including occipital, temporal and parietal cortex), right inferior frontal sulcus, and left post-central motor cortex (the familywise error rate was controlled at the voxel-level $p < 0.001$). **B)** To isolate voxels most strongly representing the stimulus features, we raised the statistical threshold, resulting in the ROI illustrated in yellow. **C)** “5/4” Dataset Binary Feature Decoding. Red dots indicate scores from individual participants. **D)** SHJ Binary Feature Decoding. The same ROI (B) was used in both datasets.

was chosen for the searchlight analysis to improve computational efficiency).³ Each feature could be decoded at rates significantly above chance (shape: $M = 0.60$, $SE = 0.02$, $t(19) = 5.78$, $p < 0.001$, size: $M = 0.70$, $SE = 0.02$, $t(19) = 9.29$, $p < 0.001$, color: $M = 0.54$, $SE = 0.01$, $t(19) = 3.64$, $p = 0.002$, position: $M = 0.91$, $SE = 0.02$, $t(19) = 25.96$, $p < 0.001$).

Next, we investigated whether the decoding accuracy of the individual perceptual dimensions covaried with the GCM attentional parameters. To do so, we fit a mixed-effects linear regression analysis (as implemented in the lme4 package for R) using restricted maximum likelihood (ReML). We included fixed-effects terms for the

³The C parameter modulates the penalty associated with training error. With large values, the classifier will choose a small-margin hyperplane, and training accuracy will be high. With smaller values, out-of-sample performance is often improved, but more training samples may be misclassified. C=1 is a common default setting for fMRI.

intercept, the attentional weights, and each visual dimension (e.g., “color”). We also included random effects terms (which were free to vary between participants) for the intercept and the attention weight parameters. This allowed us to control for baseline differences in decoding accuracy between participants, and for shared (group-level) differences in decoding accuracy between visual dimensions. We used the Kenward-Roger approximation (Kenward & Roger, 1997) to estimate degrees of freedom (reported below), and used single-sample t -tests to calculate p -values for each coefficient (using the `pbkrtest` package for R; Halekoh & Højsgaard, 2014)⁴. We computed 95% confidence intervals using bootstrap resampling (1000 simulations). The decoding accuracy of each stimulus dimension positively covaried with the behaviorally-derived GCM parameters ($b = 0.08$, 95% CI = [0.01, 0.16], $SE = 0.04$, $t(28.71) = 2.26$, $p = 0.032$), indicating that the decoding accuracy of these representations reflected their importance during decision-making.

To investigate the sensitivity of occipitotemporal feature representations to individual differences in GCM attentional weights, we conducted a permutation test. This involved shuffling the attentional weight parameters between participants (i.e., swapping the weights derived from one participant with those derived from another), and repeating the regression analysis (described above) 10,000 times. On each permutation, the correspondence for category dimensions (i.e., the dimensions depicted in Table 1, as opposed to the stimulus dimensions illustrated in Figure 1) was preserved, such that the dimensional weights derived from the behavior of one participant were assigned to the same dimensions, but to a different participant.

The unpermuted beta coefficient ($b = 0.08$) was significantly greater than those composing the null distribution ($P = 0.994$), indicating that the decoding accuracy of the occipitotemporal representations was sensitive to between-subject differences in the attentional weights. This could reflect idiosyncratic differences in behavioral strategy, and/or effects associated with perceptual saliency. Therefore, to investigate whether visual saliency may have influenced attention, we conducted a repeated measures

⁴This provides a more conservative test than the likelihood ratio test or the Wald approximation (Luke, 2016).

ANOVA for the perceptual features. There was no significant relationship between these visual features and the attentional parameters ($F(3,57) = 0.68$, $p = 0.56$). A Bayesian repeated-measures ANOVA (Rouder, Morey, Verhagen, Swagman, & Wagenmakers, 2017), additionally indicated that the null model was 4.65 times more likely than the alternative hypothesis. These results provide evidence that the observed effects were not driven by visual characteristics of the stimulus features.

SHJ Dataset. First, we confirmed that each stimulus feature could be decoded significantly above chance from the ROI illustrated in Figure 3.B. Using a four-fold, leave-one-run out cross-validation strategy, we used a linear support vector classifier ($C=1$) to decode each visual feature across all runs (including both early and late learning epochs), retaining only estimates for the last two runs (which corresponded to the late-learning phase in which behavior had largely stabilized). This four-fold cross-validation strategy yielded better decoding accuracy than a two-fold approach based on only the last two runs. This improvement reflects the increased amount of training data available in the 4-fold approach, and suggests that the multivariate patterns reflecting the individual visual features were stable across learning. Each feature could be decoded at rates significantly above chance (Figure 3.D; antennae: $M = 0.57$, $t(20) = 3.82$, $p = 0.001$; mandibles: $M = 0.56$, $t(20) = 3.22$, $p = 0.004$; legs: $M = 0.58$, $t(20) = 4.17$, $p < 0.001$).

Next, we investigated whether the decoding accuracy associated with the features covaried with SUSTAIN’s attentional parameters. To do so, we used a mixed-effects linear regression analysis to predict decoding accuracy from attention weight, visual dimension, run and rule. As described in the Methods section, distinct attentional weights were derived for each subject and each rule. The decoding accuracy for each separate run was included in the analysis. The model included fixed-effects parameters for these four variables, and random-effects parameters for the intercept, attention weight, and run (which were free to vary by participant). This allowed us to control for differences in decoding accuracy across visual dimensions and participants (as with the model used for the “5/4” dataset), while additionally controlling for effects of rule and

idiosyncratic differences in behavioral performance during the last two runs. Mirroring the findings from the “5/4” dataset, we found that the decoding accuracy of these patterns positively covaried with the attention parameters derived from SUSTAIN ($b = 0.09$, 95% CI = [0.004, 0.17], $SE = 0.04$, $t(61) = 2.13$, $p = 0.038$).

To investigate the sensitivity of occipitotemporal feature representations to individual differences in SUSTAIN’s attentional parameters, we conducted a permutation test similar to that described above (i.e., for the “5/4” experiment). This involved shuffling the attentional weight parameters between participants 10,000 times (preserving the correspondence for both rule and abstract feature). This means that the attentional weight derived from the behavior of one participant, for one particular rule and one particular category feature, was assigned to the same rule and feature, but to a different participant. The slope parameter associated with the unpermuted data ($b = 0.09$) was significantly greater than those composing the permuted null distribution ($P = 0.979$), suggesting that the visual feature representations were sensitive to idiosyncratic differences in attentional weights. A repeated measures ANOVA indicated that the perceptual dimensions did not influence the attentional parameters ($F(2,44) = 1.27$, $p = 0.291$). A Bayesian repeated measures ANOVA additionally indicated that the null model was 1.98 times more likely than the alternative hypothesis, providing evidence that the attentional weights were not influenced by visual properties of the stimulus features.

Discussion

Although differing substantially in how concepts are represented (e.g., as exemplars, prototypes, or clusters), formal categorization theories (e.g., Kruschke, 1992; Love et al., 2004; Nosofsky, 1986) tend to share a similar conception of selective attention. In these models, conceptual knowledge contorts multidimensional psychological space such that differences along behaviorally-relevant dimensions are accentuated, and differences along irrelevant dimensions are down-weighted (Figure 1, and Equations 1 & 4 in appendix B). In two datasets (Mack et al., 2016, 2013), we

evaluated the neurobiological plausibility of this idea by investigating whether occipitotemporal stimulus feature representations covaried with attention parameters derived from formal categorization models. We found that this effect was not only apparent at the group-level, but was sufficiently sensitive to reflect individual differences in conceptual knowledge.

Several previous studies have demonstrated that occipitotemporal stimulus representations are modulated by selective attention (e.g., Buffalo et al., 2010; Jehee et al., 2011; Kamitani & Tong, 2005, 2006; Luck et al., 1997; Motter, 1993; Reynolds & Chelazzi, 2004; Reynolds, Pasternak, & Desimone, 2000) and by learned conceptual knowledge (e.g., Folstein et al., 2013; Li et al., 2007; Sigala & Logothetis, 2002). These studies have relied on statistically-powerful contrastive approaches, in which representations of attended stimulus dimensions are compared to those of unattended dimensions. A general finding is that attended stimulus dimensions are more easily decoded than those that are unattended. This implies that occipitotemporal representational space might resemble that conceptualized by formal categorization theory (e.g., Kruschke, 1992; Love et al., 2004; Nosofsky, 1986). Specifically, the expansion and contraction of this space might closely reflect individual differences in the importance assigned to each stimulus dimension. However, as the contrastive approach defines selective attention with regards to the experimental paradigm, it is insensitive to individual differences in categorization strategy (e.g., Craig & Lewandowsky, 2012; Little & McDaniel, 2015; McDaniel et al., 2014; Raijmakers et al., 2014). Here, we link individual differences in behavior to individual differences in neural representation through consideration of the attentional parameters derived from formal categorization models.

We are not the first to link brain and behavior via latent model parameters. In the perceptual decision-making literature, for instance, several groups have fit the drift diffusion model (Ratcliff, 1978) to behavioral data, and identified regions of the brain where the BOLD response reflects variation in its drift rate, bias, and threshold parameters (e.g., Forstmann et al., 2008; Mulder, Wagenmakers, Ratcliff, Boekel, &

Forstmann, 2012; Purcell et al., 2010). As in the present study, several of these studies demonstrated that individual differences in behavioral strategy are reflected in the brain. Instead of linking latent model parameters to univariate BOLD amplitude, however, we used MVPA to link latent parameters to multivoxel representations of the stimulus features. This provided a precise test of the idea that selective attention contorts neural representational space.

These endogenous attentional effects are thought to arise through communication with other areas of the brain. In lateral frontal cortex, for instance, effects of endogenous attention occur earlier in time than in occipitotemporal cortex (Baldauf & Desimone, 2014; Bichot, Heard, Degennaro, & Desimone, 2015; Zhou & Desimone, 2011). Inactivation of these frontal regions (e.g., ventral prearcuate sulcus, or entire lateral prefrontal cortex) has also been associated with a reduction in the magnitude of attentional effects in occipitotemporal cortex (Bichot et al., 2015; Gregoriou, Rossi, Ungerleider, & Desimone, 2014). Interestingly, contextually-sensitive effects of endogenous attention have also been observed in the lateral geniculate nucleus (LGN), suggesting that some aspects of attention precede those in cortex (McAlonan, Cavanaugh, & Wurtz, 2008; O'Connor, Fukui, Pinsk, & Kastner, 2002; Saalmann & Kastner, 2011).

Finally, it is worth noting that, although we observed effects of selective attention across two different stimulus sets (abstract shapes in the “5/4” experiment, and insects in the SHJ experiment), and across multiple category structures (the “5/4” problem described by Medin and Schaffer (1978), and the Type I, II and VI problems described by Shepard et al. (1961)), these effects might not be apparent for all stimuli and tasks. For instance, although category training can improve perceptual discriminability of relevant stimulus features when stimuli consist of perceptually-separable features (Garner, 1976), this may not occur for integral dimensions (Op de Beeck, Wagemans, & Vogels, 2003) or for stimuli defined according to “blended” stimulus morphspaces (Folstein et al., 2013). More work is needed to better understand how attention influences occipitotemporal representations for such stimuli. One possibility is that

selective attention does not warp *perceptual* representations of integral stimulus dimensions, but might operate on abstract cognitive or “decisional” representations in higher-order cortex (Jiang et al., 2007; Nosofsky, 1987).

Conclusions

Category training is known to induce changes in both perceptual (Folstein, Gauthier, & Palmeri, 2012; Goldstone, 1994; Goldstone, Steyvers, & Larimer, 1996; Gureckis & Goldstone, 2008; Op de Beeck et al., 2003) and neural sensitivity (e.g., Dieciuc, Roque, & Folstein, 2017; Folstein et al., 2013; Folstein, Palmeri, Gauthier, & Van Gulick, 2015; Li et al., 2007; Sigala & Logothetis, 2002). In two datasets, we demonstrate that occipitotemporal stimulus representations covary with the attentional parameters derived from formal categorization theory. This effect was sufficiently sensitive to reflect individual differences in conceptual knowledge, which implies that these occipitotemporal representations are embedded within a space closely resembling that predicted by formal categorization theory (e.g., Kruschke, 1992; Love et al., 2004; Nosofsky, 1986).

By linking brain and behavior through the latent attentional parameters of cognitive models, we also link two (somewhat) disparate literatures. In the neuroscience literature, effects of selective attention are typically examined using highly-structured decision problems, and selective attention is investigated by contrasting different aspects of the experimental design (i.e., relevant vs. irrelevant stimulus dimensions). In the cognitive categorization literature, researchers have focused on developing models that accurately account for behavioral patterns of generalization across different goals and tasks. Our results indicate that these cognitive models can be used to examine effects of selective attention in the brain. This is the case, even for ill-defined decision problems (such as the “5/4” task), as the models are able to successfully account for individual differences in conceptual knowledge.

Context

Brad Love has a longstanding interests in models of categorization. He developed the SUSTAIN model (Love et al., 2004) used here, and subsequently became interested in how to theoretically relate such models to the brain (Love & Gureckis, 2007). Later, he used category category learning models in model-based fMRI analyses, such as in the two papers from which this contribution draws its data (Mack et al., 2016, 2013). Through several papers, Kurt Braunlich has investigated neurobiological mechanisms associated with categorization and generalization. Recently (Braunlich, Liu, & Seger, 2017), he found that occipitotemporal category representations are highly flexible, in that they are sensitive to transient generalization demands (i.e., strict vs. lax decision criteria). This dovetails with the present work, which examines attentional effects associated with task demands. The present work was presented at the Society for Neuroscience annual meeting (2017), and at the Cognitive Computational Neuroscience annual meeting (2017). An earlier version of the manuscript is available on the BioRxiv preprint server (Braunlich & Love, 2018).

References

- Avants, B., Tustison, N., & Johnson, H. (2009). Advanced Normalization Tools (ANTS). *Insight Journal*, 1–35.
- Baldauf, D., & Desimone, R. (2014). Neural mechanisms of object-based attention. *Science*, *344*(6182), 424–7. doi: 10.1126/science.1247003
- Bichot, N. P., Heard, M. T., Degennaro, E. M., & Desimone, R. (2015). A source for feature-based attention in the prefrontal cortex. *Neuron*, *88*, 1–13. doi: 10.1016/j.neuron.2015.10.001
- Braunlich, K., Liu, Z., & Seger, C. A. (2017). Occipitotemporal category representations are sensitive to abstract category boundaries defined by generalization demands. *The Journal of Neuroscience*, *37*(32), 3825–16.
- Braunlich, K., & Love, B. C. (2018). Occipitotemporal representations reflect individual differences in conceptual knowledge. *BioRxiv*, 264895. doi: 10.1101/264895
- Brincat, S. L., Siegel, M., Nicolai, C. V., & Miller, E. K. (2017). Gradual progression from sensory to task-related processing in cerebral cortex. doi: 10.1101/195602
- Buffalo, E. A., Fries, P., Landman, R., Liang, H., & Desimone, R. (2010). A backward progression of attentional effects in the ventral stream. *Proceedings of the National Academy of Sciences*, *107*(1), 361–365. doi: 10.1073/pnas.0907658106
- Craig, S., & Lewandowsky, S. (2012). Whichever way you choose to categorize, working memory helps you learn. *Quarterly Journal of Experimental Psychology*, *65*(3), 439–464. doi: 10.1080/17470218.2011.608854
- Dieciuc, M., Roque, N. A., & Folstein, J. R. (2017). Changing similarity: Stable and flexible modulations of psychological dimensions. *Brain Research*, *1670*, 208–219. doi: 10.1016/j.brainres.2017.06.026
- Etzel, J. A., Zacks, J. M., & Braver, T. S. (2013). Searchlight analysis: Promise, pitfalls, and potential. *NeuroImage*, *78C*, 261–269. doi: 10.1016/j.neuroimage.2013.03.041
- Folstein, J. R., Gauthier, I., & Palmeri, T. J. (2012). How category learning affects

- object representations: Not all morphspaces stretch alike. *Journal of Experimental Psychology. Learning, Memory, and Cognition*, *38*(4), 807–20. doi: 10.1037/a0025836
- Folstein, J. R., Palmeri, T. J., & Gauthier, I. (2013). Category learning increases discriminability of relevant object dimensions in visual cortex. *Cerebral Cortex*, *23*(4), 814–823. doi: 10.1093/cercor/bhs067
- Folstein, J. R., Palmeri, T. J., Gauthier, I., & Van Gulick, A. E. (2015). Category learning stretches neural representations in visual cortex. *Current Directions in Psychological Science*, *24*(1), 17–23. doi: 10.1177/0963721414550707
- Forstmann, B. U., Dutilh, G., Brown, S., Neumann, J., von Cramon, D. Y., Ridderinkhof, K. R., & Wagenmakers, E.-J. (2008). Striatum and pre-SMA facilitate decision-making under time pressure. *Proceedings of the National Academy of Sciences of the United States of America*, *105*(45), 17538–42. doi: 10.1073/pnas.0805903105
- Garner, W. (1976). Interaction of stimulus dimensions in concept and choice processes. *Cognitive Psychology*, *8*(1), 98–123. doi: 10.1016/0010-0285(76)90006-2
- Goldstone, R. L. (1994). Influences of categorization on perceptual discrimination. *Journal of Experimental Psychology. General*, *123*(2), 178–200. doi: 10.1037/0096-3445.123.2.178
- Goldstone, R. L. (2003). Learning to perceive while perceiving to learn. In R. Kimchi, M. Behrmann, & C. Olson (Eds.), *Perceptual organization in vision: Behavioral and neural perspectives*. (pp. 233–278). New Jersey: Lawrence Erlbaum Associates.
- Goldstone, R. L., Steyvers, M., & Larimer, K. (1996). Categorical perception of novel dimensions. In *Proceedings of the eighteenth annual conference of the cognitive science society* (pp. 243–248). doi: 10.1080/713756735
- Gregoriou, G. G., Rossi, A. F., Ungerleider, L. G., & Desimone, R. (2014). Lesions of prefrontal cortex reduce attentional modulation of neuronal responses and synchrony in V4. *Nature Neuroscience*, *17*(7), 1003–11. doi: 10.1038/nn.3742

- Gureckis, T. M., & Goldstone, R. L. (2008). The effect of the internal structure of categories on perception. *Proceedings of the 30th Annual Conference of the Cognitive Science Society*, 1876–1881.
- Halekoh, U., & Højsgaard, S. (2014). A Kenward-Roger approximation and parametric bootstrap methods for tests in linear mixed models - The R package pbkrtest. *Journal of Statistical Software*, *59*(9), 1–32. doi: 10.18637/jss.v059.i09
- Haxby, J. V., Gobbini, M. I., Furey, M. L., Ishai, A., Schouten, J. L., & Pietrini, P. (2001). Distributed and overlapping representations of face and objects in ventral temporal cortex. *Science*, *293*(5539), 2425–2430. doi: 10.1126/science.1063736
- Haynes, J.-D. (2015). A primer on pattern-based approaches to fMRI: Principles, pitfalls, and perspectives. *Neuron*, *87*(2), 257–270. doi: 10.1016/j.neuron.2015.05.025
- Hedge, C., Powell, G., & Sumner, P. (2017). The reliability paradox: Why robust cognitive tasks do not produce reliable individual differences. *Behavior Research Methods*, 1-21. doi: 10.3758/s13428-017-0935-1
- Jehee, J. F. M., Brady, D. K., & Tong, F. (2011). Attention improves encoding of task-relevant features in the human visual cortex. *The Journal of Neuroscience*, *31*(22), 8210–8219. doi: 10.1523/JNEUROSCI.6153-09.2011
- Jiang, X., Bradley, E., Rini, R. A., Zeffiro, T., VanMeter, J., & Riesenhuber, M. (2007). Categorization training results in shape- and category-selective human neural plasticity. *Neuron*, *53*(6), 891–903. doi: 10.1016/j.neuron.2007.02.015
- Johansen, M. K., & Palmeri, T. J. (2002). Are there representational shifts during category learning? *Cognitive Psychology*, *45*(4), 482-553. doi: 10.1016/S0010-0285(02)00505-4
- Kamitani, Y., & Tong, F. (2005). Decoding the visual and subjective contents of the human brain. *Nature Neuroscience*, *8*(5), 679–685. doi: 10.1038/nn1444
- Kamitani, Y., & Tong, F. (2006). Decoding seen and attended motion directions from activity in the human visual cortex. *Current Biology*, *16*(11), 1096–1102. doi: 10.1016/j.cub.2006.04.003

- Kenward, M. G., & Roger, J. H. (1997). Small sample inference for fixed effects from restricted maximum likelihood. *Biometrics*, *53*(3), 983–997.
- Kriegeskorte, N., Goebel, R., & Bandettini, P. (2006). Information-based functional brain mapping. *Proceedings of the National Academy of Sciences of the United States of America*, *103*(10), 3863–8. doi: 10.1073/pnas.0600244103
- Kriegeskorte, N., & Kievit, R. A. (2013). Representational geometry: Integrating cognition, computation, and the brain. *Trends in cognitive sciences*, *17*(8), 401–12. doi: 10.1016/j.tics.2013.06.007
- Kriegeskorte, N., Mur, M., & Bandettini, P. A. (2008). Representational similarity analysis: Connecting the branches of systems neuroscience. *Frontiers in Systems Neuroscience*, *2*. doi: 10.3389/neuro.06.004.2008
- Kruschke, J. (1992). ALCOVE: An exemplar-based connectionist model of category learning. *Psychological Review*.
- Li, S., Ostwald, D., Giese, M., & Kourtzi, Z. (2007). Flexible coding for categorical decisions in the human brain. *The Journal of Neuroscience*, *27*(45), 12321–30. doi: 10.1523/JNEUROSCI.3795-07.2007
- Little, J. L., & McDaniel, M. (2015). Individual differences in category learning: Memorization versus rule abstraction. *Memory & cognition*, *43*(2), 283–97. doi: 10.3758/s13421-014-0475-1
- Love, B. C., & Gureckis, T. M. (2007). Models in search of a brain. *Cognitive, Affective & Behavioral Neuroscience*, *7*(2), 90–108. doi: 10.3758/CABN.7.2.90
- Love, B. C., Medin, D. L., & Gureckis, T. M. (2004). SUSTAIN: A network model of category learning. *Psychological Review*, *111*(2), 309–32. doi: 10.1037/0033-295X.111.2.309
- Luck, S. J., Chelazzi, L., Hillyard, S. A., & Desimone, R. (1997). Neural mechanisms of spatial selective attention in areas V1, V2, and V4 of macaque visual cortex. *Journal of Neurophysiology*, *77*(1), 24–42.
- Luke, S. G. (2016). Evaluating significance in linear mixed-effects models in R. *Behavior Research Methods*, *49*(4), 1494–1502.

- Mack, M. L., Love, B. C., & Preston, A. R. (2016). Dynamic updating of hippocampal object representations reflects new conceptual knowledge. *Proceedings of the National Academy of Sciences of the United States of America*, *113*(46), 13203–13208. doi: 10.1073/pnas.1614048113
- Mack, M. L., Preston, A. R., & Love, B. C. (2013). Decoding the brain's algorithm for categorization from its neural implementation. *Current Biology*, *23*(20), 2023–2027. doi: 10.1016/j.cub.2013.08.035
- McAlonan, K., Cavanaugh, J., & Wurtz, R. H. (2008). Guarding the gateway to cortex with attention in visual thalamus. *Nature*, *456*(7220), 391–394. doi: 10.1038/nature07382
- McDaniel, M. A., Cahill, M. J., Robbins, M., & Wiener, C. (2014). Individual differences in learning and transfer: Stable tendencies for learning exemplars versus abstracting rules. *Journal of Experimental Psychology.*, *143*, *143*(2, 2), 668, 668–693. doi: 10.1037/a0032963, 10.1037/a0032963
- Medin, D. L., & Schaffer, M. M. (1978). Context theory of classification learning. *Psychological Review*, *85*(3), 207–238. doi: 10.1037/0033-295X.85.3.207
- Miller, E. K., & Cohen, J. D. (2001). An integrative theory of prefrontal cortex function. *Annual Review of Neuroscience*, *24*, 167–202. doi: 10.1146/annurev.neuro.24.1.167
- Minda, J. P., & Smith, J. D. (2002). Comparing prototype-based and exemplar-based accounts of category learning and attentional allocation. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *28*(2), 275–292. doi: 10.1037//0278-7393.28.2.275
- Misaki, M., Kim, Y., Bandettini, P. A., & Kriegeskorte, N. (2010). Comparison of multivariate classifiers and response normalizations for pattern-information fMRI. *NeuroImage*, *53*(1), 103–118. doi: 10.1016/j.neuroimage.2010.05.051
- Motter, B. C. (1993). Focal attention produces spatially selective processing in visual cortical areas V1, V2, and V4 in the presence of competing stimuli. *Journal of Neurophysiology*, *70*(3), 909–919. doi: 0022-3077/93

- Mulder, M. J., Wagenmakers, E.-J., Ratcliff, R., Boekel, W., & Forstmann, B. U. (2012). Bias in the brain: A diffusion model analysis of prior probability and potential payoff. *Journal of Neuroscience*, *32*(7), 2335–2343. doi: 10.1523/JNEUROSCI.4156-11.2012
- Mumford, J. A., Turner, B. O., Ashby, F. G., & Poldrack, R. A. (2012). Deconvolving BOLD activation in event-related designs for multivoxel pattern classification analyses. *NeuroImage*, *59*(3), 2636–43. doi: 10.1016/j.neuroimage.2011.08.076
- Nosofsky, R. M. (1986). Attention, similarity, and the identification-categorization relationship. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *115*(1), 39–57. doi: 10.1037/0278-7393.13.1.87
- Nosofsky, R. M. (1987). Attention and learning processes in the identification and categorization of integral stimuli. *Journal of Experimental Psychology: Learning, Memory & Cognition*, *13*(1), 87–108. doi: 10.1037/0278-7393.13.1.87
- Nosofsky, R. M., & Zaki, S. R. (2002). Exemplar and prototype models revisited: Response strategies, selective attention, and stimulus generalization. *Journal of Experimental Psychology-Learning Memory and Cognition*, *28*(5), 924–940. doi: 10.1037/0278-7393.28.5.924
- O'Bryan, S. R., Walden, E., Serra, M. J., & Davis, T. (2018). Rule activation and ventromedial prefrontal engagement support accurate stopping in self-paced learning. *NeuroImage*, *172*, 415–426. doi: 10.1016/j.neuroimage.2018.01.084
- O'Connor, D. H., Fukui, M. M., Pinsk, M. A., & Kastner, S. (2002). Attention modulates responses in the human lateral geniculate nucleus. *Nature Neuroscience*, *5*(11), 1203–1209. doi: 10.1038/nn957
- Op de Beeck, H., Wagemans, J., & Vogels, R. (2003). The effect of category learning on the representation of shape: Dimensions can be biased but not differentiated. *Journal of Experimental Psychology. General*, *132*(4), 491–511. doi: 10.1037/0096-3445.132.4.491
- Palmeri, T. J., Love, B. C., & Turner, B. M. (2017). Model-based cognitive neuroscience. *Journal of Mathematical Psychology*, *76*, 59–64. doi:

10.1016/j.jmp.2016.10.010

- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., . . . Vanderplas, J. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, *12*, 2825–2830.
- Purcell, B. A., Heitz, R. P., Cohen, J. Y., Schall, J. D., Logan, G. D., & Palmeri, T. J. (2010). Neurally constrained modeling of perceptual decision making. *Psychological Review*, *117*(4), 1113–1143.
- Raijmakers, M. E. J., Schmittmann, V. D., & Visser, I. (2014). Costs and benefits of automatization in category learning of ill-defined rules. *Cognitive Psychology*, *69*, 1–24. doi: 10.1016/j.cogpsych.2013.12.002
- Ratcliff, R. (1978). A theory of memory retrieval. *Psychological Review*, *85*(2), 59:108.
- Rehder, B., & Hoffman, A. B. (2005a). Eyetracking and selective attention in category learning. *Cognitive Psychology*, *51*(1), 1–41. doi: 10.1016/j.cogpsych.2004.11.001
- Rehder, B., & Hoffman, A. B. (2005b). Thirty-something categorization results explained: Selective attention, eyetracking, and models of category learning. *Journal of Experimental Psychology: Learning, Memory & Cognition*, *31*(5), 811–29. doi: 10.1037/0278-7393.31.5.811
- Reynolds, J. H., & Chelazzi, L. (2004). Attentional modulation of visual processing. *Annual Review of Neuroscience*, *27*(1), 611–647. doi: 10.1146/annurev.neuro.26.041002.131039
- Reynolds, J. H., Pasternak, T., & Desimone, R. (2000). Attention increases sensitivity of V4 neurons. *Neuron*, *26*(3), 703–714. doi: 10.1016/S0896-6273(00)81206-4
- Rouder, J. N., Morey, R. D., Verhagen, J., Swagman, A. R., & Wagenmakers, E.-J. (2017). Bayesian analysis of factorial designs. *Psychological Methods*, *22*(2), 304–321. doi: 10.1037/met0000057
- Saalmann, Y. B., & Kastner, S. (2011). Cognitive and perceptual functions of the visual thalamus. *Neuron*, *71*.
- Shepard, R. N. (1957). Stimulus and response generalization: A stochastic model relating generalization to distance in psychological space. *Psychometrika*, *22*(4),

- 325–345. doi: 10.1007/BF02288967
- Shepard, R. N., Hovland, C., & Jenkins, H. (1961). Learning and memorization of classifications. *Psychological Monographs: General and Applied*, *75*(13).
- Sigala, N., & Logothetis, N. K. (2002). Visual categorization shapes feature selectivity in the primate temporal cortex. *Nature*, *415*(6869), 318–320. doi: 10.1038/415318a
- Smith, J. D., & Minda, J. P. (1998). Prototypes in the mist: The early epochs of category learning. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *24*(6), 1411–1436. doi: 10.1037//0278-7393.24.6.1411
- Townsend, J. T., & Ashby, F. G. (1982). Experimental test of contemporary mathematical models of visual letter recognition. *Journal of Experimental Psychology. Human Perception and Performance*, *8*(6), 834–854. doi: 10.1037/0096-1523.8.6.834
- Tsotsos, J. K. (2011). *A Computational Perspective on Visual Attention*. MIT Press.
- Turner, B. M., Forstmann, B. U., Love, B. C., Palmeri, T. J., & Van Maanen, L. (2017). Approaches to analysis in model-based cognitive neuroscience. *Journal of Mathematical Psychology*, *76*, 65–79. doi: 10.1016/j.jmp.2016.01.001
- Tversky, A. (1977). Features of similarity. *Psychological Review*, *84*(4), 327:352.
- Winkler, A. M., Ridgway, G. R., Webster, M. A., Smith, S. M., & Nichols, T. E. (2014). Permutation inference for the general linear model. *NeuroImage*, *92*, 381–397. doi: 10.1016/j.neuroimage.2014.01.060
- Zaki, S. R., Nosofsky, R. M., Stanton, R. D., & Cohen, A. L. (2003). Prototype and exemplar accounts of category learning and attentional allocation: A reassessment. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *29*(6), 1160–1173. doi: 10.1037/0278-7393.29.6.1160
- Zhou, H., & Desimone, R. (2011). Feature-Based attention in the frontal eye field and area V4 during visual search. *Neuron*, *70*(6), 1205-1217. doi: 10.1016/j.neuron.2011.04.032

Appendix A: Searchlight Results

Size	x	y	z	t	BA	Region
23972	14	-74	4	12.717	18	Calcarine_R
	28	-62	54	10.5	7	Parietal_Sup_R
	-54	-18	44	7.256	3	Postcentral_L
233	50	24	26	7.447	48	Frontal_Inf_Tri_R
22	-46	12	44	6.621	9	Frontal_Mid_L

Table 3

“5/4” Dataset Binary Feature Decoding: Searchlight Results. The familywise error rate was controlled at the voxel level ($p < 0.001$)

Appendix B: Computational Models and Attentional Parameters

For the first dataset (Mack et al., 2013), we considered the Generalized Context Model (GCM; Nosofsky, 1986), which posits that conceptual knowledge consists of memory for individual exemplars. For the SHJ experiment (Mack et al., 2016), we considered the attentional parameters from SUSTAIN (Supervised and Unsupervised STRatified Adaptive Incremental Network; Love et al., 2004). Details about the models can be found in the original papers. Here, we provide a brief overview of each.

GCM

In the Generalized Context Model (GCM; Nosofsky, 1987), the psychological distance, d between stimuli i and j can be calculated as the attentionally-weighted sum of their unsigned differences across dimensions, k :

$$d_{ij} = \sum_k [w_k |x_{ik} - x_{jk}|^r]^{1/r}, \quad (1)$$

where w indicate the attentional parameters assigned to each dimension. The r parameter is set to 1 (city-block distance) for perceptually separable stimulus dimensions (as in the “5/4” dataset), and r is set to 2 (Euclidean distance) for integral dimensions (Garner, 1976). Similarity is an exponentially-decaying function of psychological distance:

$$s_{ij} = d^{-cd_{ij}}, \quad (2)$$

where the shape of the similarity gradient is influenced by the sensitivity parameter, c . The probability of choosing category “A”, given stimulus, i , is given by the choice rule:

$$P(A|i) = \frac{\left(\sum_{a \in A} s_{ia}\right)^\gamma}{\left(\sum_{a \in A} s_{ia}\right)^\gamma + \left(\sum_{b \in B} s_{ib}\right)^\gamma}, \quad (3)$$

where γ governs the degree of deterministic responding.

SUSTAIN

SUSTAIN is a semi-supervised clustering model, which incrementally learns to solve categorization problems by first applying simple solutions, and then increasing complexity as required. Through experience, the model can learn to group similar items into common clusters, and can make inferences about novel stimuli based on its perceptual similarity to existing clusters (i.e., based on perceptual similarity, clusters compete to predict latent stimulus attributes). When unexpected feedback is received, the model can also learn in a supervised fashion by creating a new cluster to represent the novel stimulus.

In SUSTAIN, all clusters contain receptive fields (RF's) for each stimulus dimension. As new stimuli are added to the cluster, the model learns by adjusting the position of each RF to best match the cluster's expectation for novel stimuli. As the RF is an exponential function, a cluster's activation, α , decreases exponentially with distance from its preferred value:

$$\alpha(\mu) = \lambda e^{-\lambda\mu}, \quad (4)$$

where μ represents the distance of the stimulus dimension value from the cluster's preferred stimulus dimension value, and where λ represents the tuning (or width) of the RF. The λ parameters are specific to dimensions, but are shared across dimensions, and so, like the attentional parameters in the GCM, the λ parameters in SUSTAIN modulate the influence of each stimulus dimension on the overall decision outcome.

The overall activation of a cluster, H , involves consideration of each dimension, k :

$$H = \frac{\sum_k (\lambda_k)^\gamma e^{-\lambda_k \mu_k}}{\sum_k (\lambda_k)^\gamma}, \quad (5)$$

where the γ parameter (which is always non-negative) modulates the influence of the λ parameters on the choice outcome. When γ is large, attended dimensions (which are associated with large λ values, and narrow RF's), dominate the activation function (eq. 5); when γ is zero, the λ parameters are ignored, and all dimensions exert equal

influence on the choice.

SUSTAIN was fit to the SHJ dataset in a supervised fashion, using the same trial order experienced by the participants; it was also fit across rule-switches, such that learning from one task was carried over to the next. Thus, SUSTAIN was capable of reflecting learning, as well as carry-over effects associated with previously learned rules.