

Dynamic categorization rules alter representations in human visual cortex

Authors: Margaret M. Henderson^{1,2,3}, John T. Serences^{3,4,5}, Nuttida Rungratsameetaweemana^{3,6,7}

¹ Neuroscience Institute, Carnegie Mellon University, Pittsburgh, USA

² Department of Machine Learning, Carnegie Mellon University, Pittsburgh, USA

³ Neurosciences Graduate Program, University of California, San Diego, La Jolla, USA

⁴ Department of Psychology, University of California, San Diego, La Jolla, USA

⁵ Kavli Foundation for the Brain and Mind, University of California, San Diego, La Jolla, USA

⁶ The Salk Institute for Biological Studies, La Jolla, USA

⁷ Department of Biomedical Engineering, Columbia University, New York, USA

Abstract

Everyday perceptual tasks require sensory stimuli to be dynamically encoded and analyzed according to changing behavioral goals. For example, when searching for an apple at the supermarket, one might first find the Granny Smith apples by separating all visible apples into the categories “green” and “non-green”. However, suddenly remembering that your family actually likes Fuji apples would necessitate reconfiguring the boundary to separate “red” from “red-yellow” objects. This flexible processing enables identical sensory stimuli to elicit varied behaviors based on the current task context. While this phenomenon is ubiquitous in nature, little is known about the neural mechanisms that underlie such flexible computation.

Traditionally, sensory regions have been viewed as mainly devoted to processing inputs, with limited involvement in adapting to varying task contexts. However, from the standpoint of efficient computation, it is plausible that sensory regions integrate inputs with current task goals, facilitating more effective information relay to higher-level cortical areas. Here we test this possibility by asking human participants to visually categorize novel shape stimuli based on different linear and non-linear boundaries. Using fMRI and multivariate analyses of retinotopically-defined visual areas, we found that shape representations in visual cortex became more distinct across relevant decision boundaries in a context-dependent manner, with the largest changes in discriminability observed for stimuli near the decision boundary. Importantly, these context-driven modulations were associated with improved categorization performance. Together, these findings demonstrate that codes in visual cortex are adaptively modulated to optimize object separability based on currently relevant decision boundaries.

Keywords: context-dependent processing, decision making, human visual cortex, decision boundaries, task modulations, neural mechanisms

Introduction

Perceptual categorization is a fundamental cognitive ability that allows us to organize and understand the myriad stimuli encountered in our sensory environment. By forming categories, observers are able to generalize existing knowledge to new incoming inputs, facilitating efficient perception and decision-making (Bruner, 1957; Freedman & Assad, 2016). Within the visual system, categories can capture divisions within the natural structure of a stimulus space (Rosch et al., 1976) or can reflect the learning of arbitrary discrete boundaries along stimulus dimensions that would otherwise be represented continuously (Ashby & Maddox, 2005). At the same time, categorization in the real world is a highly dynamic cognitive process, in which the category membership of stimuli may change over time. For example, when making a categorical decision about produce at the farmer's market, depending on our goals we might think of carrots in the same category as lettuce (vegetables) or the same category as tangerines (orange colored items). Perceptual categorization is thus also tightly connected with flexible prioritization of information based on current task demands (Biederman et al., 1973; McAdams & Maunsell, 1999). Within contexts where task goals change dynamically over time, the neural mechanisms supporting categorization of sensory stimuli are not yet understood.

Past work has provided some insight into how category learning impacts representations of sensory stimuli. Behaviorally, learning to categorize stimuli in a continuous feature space can lead to perceptual changes such as an increase in sensitivity to changes along a relevant stimulus dimension, and an increase in perceptual discriminability of stimuli belonging to different categories (Goldstone, 1994; Livingston et al., 1998; Newell & Bühlhoff, 2002). Such changes are also reflected in the brain – electrophysiology studies in macaques have demonstrated that after learning of a categorization task, neurons in inferotemporal cortex (ITC) become more strongly selective for diagnostic dimensions of stimuli (Sigala & Logothetis, 2002), and neural populations in ITC also contain information encoding the learned category status of stimuli (Meyers et al., 2008; Tanaka, 1996). In human functional magnetic resonance imaging (fMRI) studies, learning to discriminate object categories has been shown to increase neural responses to objects in extrastriate cortex (Gauthier et al., 2000; Op de Beeck et al., 2006) and lead to sharpening of visual representations as measured with fMRI adaptation (Folstein et al., 2015; Folstein et al., 2013; Jiang et al., 2007). Moreover, recent work has shown that learning a decision boundary can alter representations of orientation in early visual areas, with representations becoming biased away from the decision boundary (Ester et al., 2020). At the same time, other work has suggested that the effects of category status on sensory representations are more prominent in prefrontal cortex (PFC) than visual areas. This suggests that the primary role of visual areas may be restricted to perceptual analysis, rather than decision-related processing (Freedman et al., 2003; McKee et al., 2014; Meyers et al., 2008).

From an efficient processing perspective, it is plausible that visual areas play a more active role in decision-making, potentially encoding decision-related variables, task contexts, choices, or motor outcomes. Such coding would enable visual areas to process sensory inputs in a manner conducive to downstream processing. Emerging evidence from rodent studies supports this view. For instance, activity that was thought to reflect random fluctuations in neural

representations within sensory areas has been linked to choice-related motor activities and decision outcomes (Musall et al., 2019; Stringer et al., 2019). Furthermore, recent findings indicate that early sensory areas robustly encode task context variables, such as expectations and decision rules, during dynamic decision-making tasks (Ebrahimi et al., 2022; Findling et al., 2023). Yet, the extent to which human sensory areas similarly code for task-related variables and adapt their representations contextually is unclear.

In addition, the mechanisms by which categorical decision-making flexibly shapes neural representations, particularly in tasks necessitating the switching between distinct decision rules, are not well understood. Prior work has demonstrated that neural populations in PFC can dynamically encode different boundaries depending on the currently relevant task rule (Cromer et al., 2010; Roy et al., 2010), providing one potential neural mechanism for dynamic decision-making. Similarly, a human neuroimaging study using novel objects suggested that representations in frontoparietal areas can encode different category distinctions between objects depending on their task relevance (Jackson et al., 2017). This study also found evidence for similar (albeit weaker) effects in the lateral occipital complex (LOC), suggesting that representations in visual areas may also be modified by task-relevance. Thus, it remains an open question whether and how varying task contexts interact with representations in visual cortex, as well as how these modulations may contribute to downstream task performance.

Here we address these gaps by investigating how neural responses in human visual cortex flexibly adapt to dynamic task contexts, as induced by varying categorization rules. We hypothesized that task context modulates sensory representations such that changes in the decision boundary are actively integrated during the early analysis of sensory information. To examine the effects of categorization within an abstract stimulus space, we generated a two-dimensional space of shape stimuli (Op de Beeck et al., 2001; Zahn & Roskies, 1972) that were viewed by human participants undergoing fMRI scanning. Participants categorized shapes according to different rules: linear boundaries (*Linear-1* and *Linear-2* tasks) or a non-linear boundary (*Nonlinear* task). These task contexts were interleaved across scanning runs, necessitating real-time cognitive adaptation to distinct categorization requirements applied to physically identical stimuli. Each task incorporated both "easy" and "hard" trials drawn from distinct locations in the shape space, enabling us to concurrently examine the influence of perceptual difficulty on decision processes. Using multivariate decoding methods in retinotopically-defined visual areas, we measured shape representations in each categorization task and examined how representations differed across task contexts. We predicted that shape representations would be more discriminable across a given decision boundary when that boundary was relevant for the current task. Findings from our neural data are in line with this account. Importantly, we further show that an increase in neural discriminability is directly linked to improved task performance.

Results

We trained 7 human participants to perform a shape categorization task while in the fMRI scanner, with each subject participating in 3 scanning sessions that each lasted 2 hours (Figure 1A). Shape stimuli varied parametrically along two independent axes, generating a two-dimensional shape space, and each condition of the task required shapes to be categorized according to either a linear boundary (*Linear-1* and *Linear-2* tasks) or a nonlinear boundary that required grouping together of non-adjacent quadrants (*Nonlinear* task). These different categorization tasks were performed during different scanning runs within each session, meaning that participants needed to flexibly apply different decision rules depending on the task condition for the current run (see *Methods*). Each task included a mixture of “easy” trials and “hard” trials. On the “easy” trials, a common set of 16 shapes, making up a 4x4 grid which we refer to as the main grid (black dots in Figure 1B), were shown in all tasks, while on “hard” trials, shapes were sampled from portions of the shape space near the active boundary, which made the current task more challenging (light gray dots in Figure 1B).

To verify the two-dimensional structure of our shape space, we used an image similarity analysis based on GIST features (Oliva & Torralba, 2001; see *Methods*) to assess the perceptual similarity between shape stimuli. As expected, a principal components analysis (PCA) performed on the GIST features revealed a two-dimensional grid structure, with the two shape space axes oriented roughly orthogonal to one another in PC space (Figure 1C). In addition, measuring the linear separability (based on between-category versus within-category Euclidean distances; see *Methods*) of shapes across each category boundary based on GIST features revealed that shapes were most separable across the *Linear-2* boundary, followed by the *Linear-1* boundary, with lowest separability for the *Nonlinear* boundary (Figure 1D). A similar pattern was found when computing separability using features from a self-supervised deep neural network model (SimCLR; T. Chen et al., 2020; see *Methods*), suggesting that these relationships held even when considering a broader set of image features (Figure 1D; darker purple bars). The low separability of the *Nonlinear* categories relative to the *Linear-1* and *Linear-2* categories is consistent with the *Nonlinear* boundary being nonlinear in shape space.

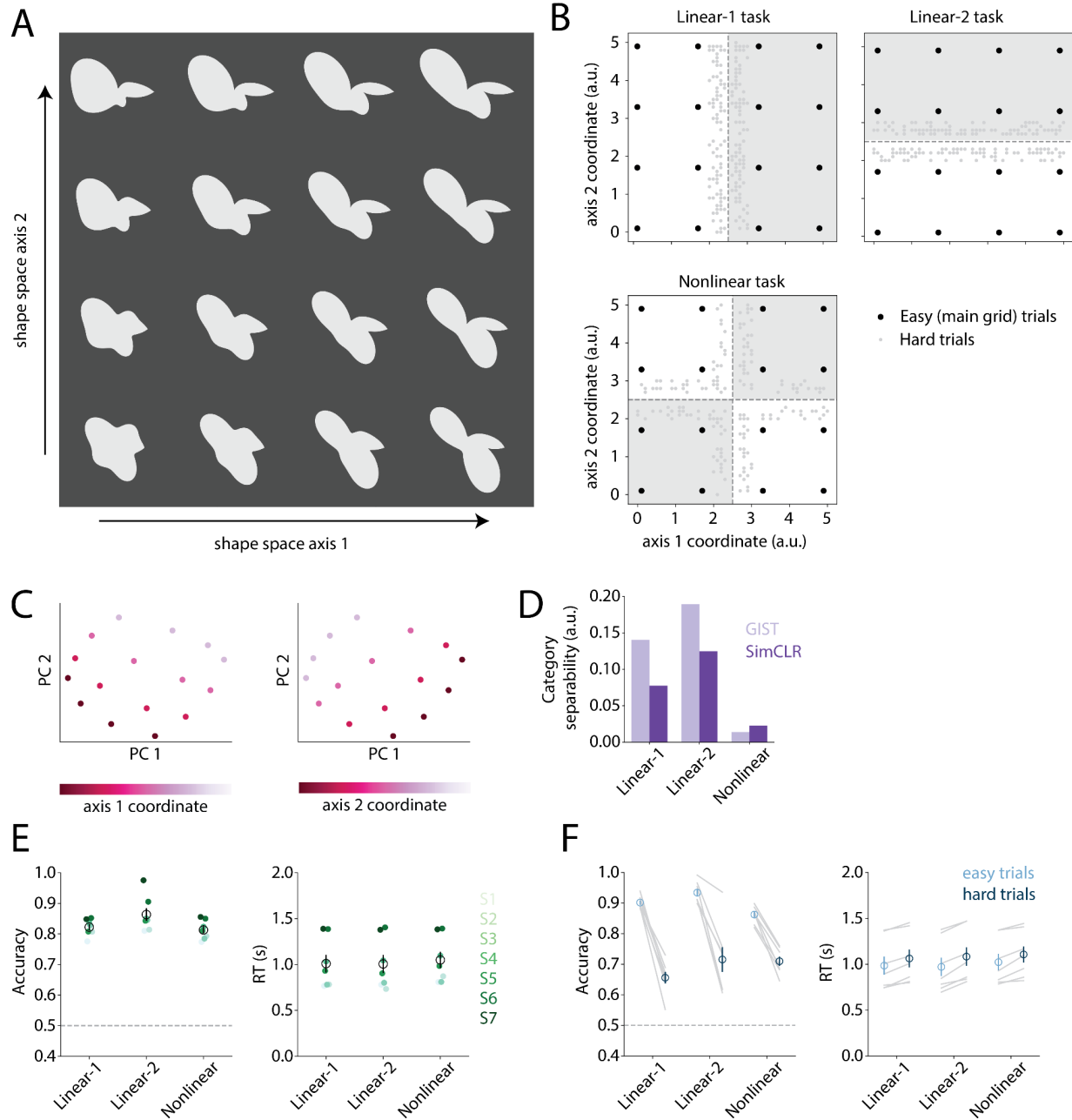


Figure 1. Stimulus set, task design, and behavioral performance. **(A)** Two-dimensional shape space used for categorization tasks in this experiment. Shapes are generated using radial frequency contours (Op de Beeck et al., 2001; Zahn & Roskies, 1972) that vary along two independent dimensions, referred to as axis 1 and axis 2. See *Methods* for more details. **(B)** Illustration of the tasks (*Linear-1*, *Linear-2*, *Nonlinear*) performed by participants while in the fMRI scanner. Points in each plot indicate the positions in shape space that were sampled, and dotted lines indicate the relevant categorization boundaries for each task. Black dots represent the 16 positions in the “main grid”, which were sampled on “easy” trials in every task, while light gray dots represent positions that were sampled on “hard” trials, which differed depending on the task. Hard trial shape positions were sampled from the region nearest the relevant

categorization boundary. Different tasks were performed during different scan runs. In each task, every trial consisted of the presentation of a single shape (1s), and participants were instructed to respond with a button press indicating which category the presented shape fell into. See *Methods* for more details on task design. **(C-D)** Image similarity analysis: we computed activations from two computer vision models, GIST (Oliva & Torralba, 2001) and SimCLR (T. Chen et al., 2020) for each of the 16 main grid shape images. **(C)** Visualization of a principal components analysis (PCA) performed on the GIST model features, where each plotted point represents one shape in PC space, colored according to the coordinate value along axis 1 (left) or axis 2 (right). **(D)** Quantification of the separability of shape categories within each feature space, computed based on the ratio of between-category to within-category Euclidean distance values. See *Methods* for more details. **(E)** Behavioral accuracy (left) and response time (RT; right) in each task. Dots in shades of green represent individual participants; open circles and error bars represent the mean \pm SEM across 7 participants. **(F)** Accuracy (left) and RT (right) for each task separated into “easy” and “hard” trials, where easy refers to trials sampling the 16 shapes in the main grid (black dots in B), and hard refers to trials sampling more challenging portions of the shape space for each task (light gray dots in B). Gray lines represent individual participants, open circles and error bars represent the mean \pm SEM across 7 participants.

Across participants, behavioral accuracy (Figure 1E) was highest for the *Linear-2* task (0.86 ± 0.02 ; mean \pm SEM across 7 participants), followed by the *Linear-1* task (0.82 ± 0.01) and the *Nonlinear* task (0.81 ± 0.01). A repeated measures ANOVA revealed a main effect of task ($F_{(2,12)} = 9.67$, $p < 0.001$; p-values obtained using permutation test; see *Methods*), and post-hoc tests showed that accuracy was significantly higher for the *Linear-2* task than all other tasks (*Linear-1* vs. *Linear-2*: $t_{(6)} = -2.65$, $p = 0.016$; *Linear-2* vs. *Nonlinear*: $t_{(6)} = 3.94$, $p = 0.018$; paired t-tests with permutation; see *Methods*). This advantage for the *Linear-2* task is consistent with the high relative separability across the *Linear-2* boundary based on image features shown in the previous analysis (Figure 1D). In terms of response times (RTs), a significant main effect of task was also found ($F_{(2,12)} = 5.28$, $p = 0.013$). No difference in RTs between the *Linear-1* and *Linear-2* tasks was observed, but RTs were significantly slower for the *Nonlinear* task than the *Linear-1* task ($t_{(6)} = -3.15$, $p = 0.050$). In addition to these differences across tasks, we also observed a consistent difference between performance on the easy and hard trials within each task (Figure 1F), which was expected based on the task design. Accuracy was significantly higher on easy versus hard trials within each task (*Linear 1*: $t_{(6)} = 11.88$, $p = 0.018$; *Linear-2*: $t_{(6)} = 5.87$, $p = 0.017$; *Nonlinear*: $t_{(6)} = 11.10$, $p = 0.016$), and RT was significantly faster on easy versus hard trials within each task (*Linear 1*: $t_{(6)} = -6.15$, $p = 0.019$; *Linear-2*: $t_{(6)} = -8.08$, $p = 0.015$; *Nonlinear*: $t_{(6)} = -4.35$, $p = 0.017$).

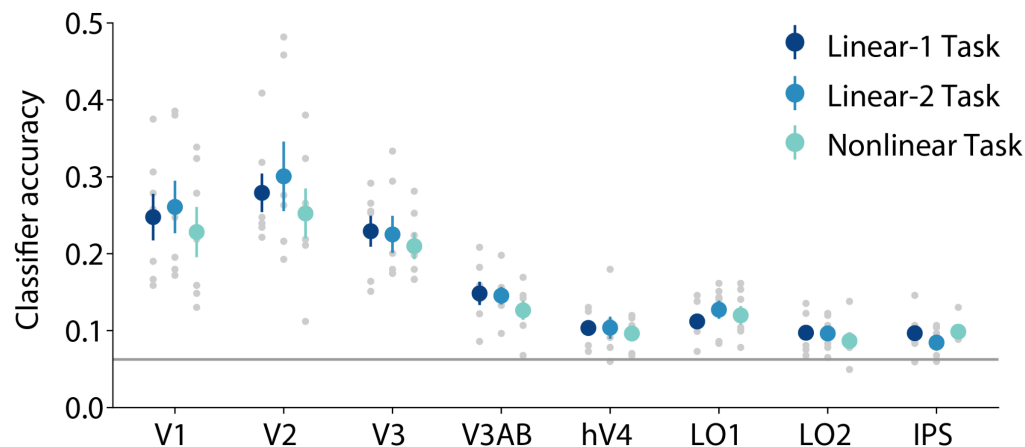


Figure 2. Overall classification accuracy. A multinomial (16-way) logistic regression classifier was trained to classify the shape shown on each trial within each task condition (*Linear-1*, *Linear-2*, and *Nonlinear* tasks). Classifiers were trained and tested within each task condition separately, training using data from the main grid trials only (i.e. black dots in Figure 1B). Plotted values reflect overall 16-way prediction accuracy of classifiers for each task and each ROI, computed using trials from the main grid only. Gray dots represent individual participants, colored circles and error bars represent the mean \pm SEM across 7 participants, horizontal line indicates chance decoding accuracy of 1/16. All classification accuracy values were above chance at the participant-averaged level (FDR corrected, $q < 0.01$); see *Methods* for more details.

Next, we examined the neural representations of shape stimuli in each task, under the hypothesis that shape representations would differ across task conditions in accordance with the changing decision boundary. To achieve this we used multivariate classification to analyze single-trial voxel activation patterns from retinotopically defined ROIs (Figure 2). We trained a 16-way multinomial classifier using L2 regularization and data from each task separately, and found that overall classification accuracy was highest in V2 (16-way accuracy averaged across tasks: 0.28 ± 0.03 ; mean \pm SEM across 7 participants), followed by V1 (0.25 ± 0.03) and V3 (0.22 ± 0.02). Participant-averaged classification accuracy was significantly above chance for every ROI in every task (significance evaluated using a permutation test; FDR corrected; all $q < 0.01$; see *Methods*).

To characterize the neural shape space, we used the output of the 16-way classifier to compute a confusion matrix for each ROI and for each task, which captures how often the classifier assigned each shape label to each shape in the test dataset (Figure 3; see *Methods*). For V1, this confusion matrix revealed that shape confusability was related to distance in shape space, with the classifier tending to make more errors between shapes that were adjacent in shape space (off-diagonal structure in Figure 3A). This relationship with distance can also be

seen by plotting the proportion of predictions as a function of the distance between predicted and actual shape space coordinates (Figure 3B). Importantly, the distances between shape space points were not specified in the construction of the classifier, where all 16 points were treated as independent categories. Thus, the emergence of this structure in the classifier confusion matrix provides evidence for a two-dimensional representation of the shape space grid in V1. A similar pattern was seen in all other ROIs tested.

Next, we examined how well the neural shape space measured in each task aligned with each decision rule. To examine this, we first constructed “template” confusion matrices for the *Linear-1* and *Linear-2* boundaries, where each template had 1 for shape pairs that were on the same side of the category boundary for that task and 0 for shape pairs that were on different sides (Figure 3C). We then correlated these template matrices with the real confusion matrices for each task (Figure 3D). This analysis revealed that the similarity of confusion matrices to each template differed depending on task. A three-way repeated measures ANOVA on the template similarity values showed main effects of ROI and Template, as well as a significant ROI x Template interaction and a significant Task x Template interaction (ROI: $F_{(7,42)} = 49.14$, $p < 0.001$; Task: $F_{(1,6)} = 5.53$, $p = 0.059$; Template: $F_{(1,6)} = 22.9$, $p = 0.003$; ROI x Task: $F_{(7,42)} = 1.48$, $p = 0.204$; ROI x Template: $F_{(7,42)} = 3.83$, $p = 0.004$; Task x Template: $F_{(1,6)} = 6.77$, $p = 0.039$; ROI x Task x Template: $F_{(7,42)} = 0.89$, $p = 0.515$; p-values obtained using permutation test; see *Methods*). Evaluating the similarity values for each template separately, we found that across all ROIs, the *Linear-2* template was significantly more similar to confusion matrices computed from the *Linear-2* task versus the *Linear-1* task (two-way repeated measures ANOVA; ROI: $F_{(7,42)} = 31.92$, $p < 0.001$; Task: $F_{(1,6)} = 9.62$, $p = 0.018$; ROI x Task: $F_{(7,42)} = 1.17$, $p = 0.352$). Post-hoc tests showed that the difference in similarity to the *Linear-2* template between the *Linear-2* and *Linear-1* tasks was significant in LO1 ($t_{(6)} = -3.41$, $p = 0.014$; paired t-test with permutation; see *Methods*). These findings suggest that shape representations in LO1 were more aligned with the *Linear-2* template when the *Linear-2* boundary was relevant than when it was irrelevant for the present task. However, the similarity of confusion matrices to the *Linear-1* template did not differ significantly across tasks (two-way repeated measures ANOVA; ROI: $F_{(7,42)} = 34.17$, $p < 0.001$; Task: $F_{(1,6)} = 0.19$, $p = 0.676$; ROI x Task: $F_{(7,42)} = 1.24$, $p = 0.303$). Additionally, when we constructed a template for the *Nonlinear* task, we did not observe a difference in the similarity of confusion matrices to the *Nonlinear* template across tasks (Supplementary Figure 1). Together, these results suggest that shape representations in visual cortex during our task may reorganize in a way that reflects the current decision boundary and shifting cognitive demands.

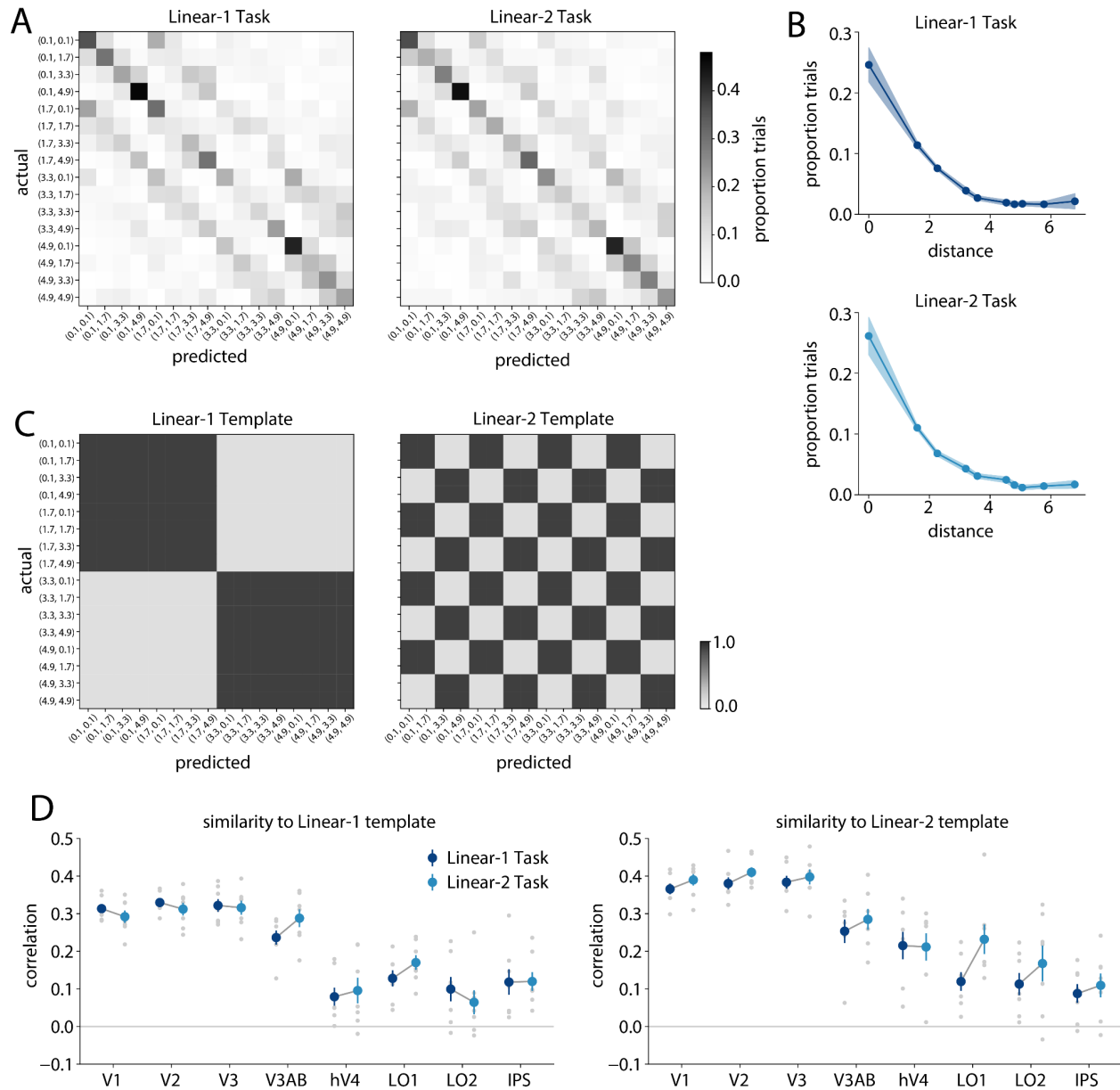


Figure 3. Classifier confusion matrices suggest restructuring of shape representations between the *Linear-1* and *Linear-2* tasks. **(A)** Classifier confusion matrices for V1 in each task, where each row represents the set of trials on which a given shape was actually shown, and the columns represent the proportion of those trials that the classifier predicted as having each of the 16 shape labels (each row sums to 1). Confusion matrices were computed using main grid trials only, and are averaged across 7 participants. **(B)** A simplified view of the classifier confusion data for V1: we computed the proportion of trials on which the actual and predicted shapes were separated by a given distance in shape space. Colored lines and shaded error bars indicate mean \pm SEM across 7 participants. **(C)** Template matrices for the *Linear-1* and *Linear-2* tasks, representing the pattern of confusability expected for a perfect binary representation of each categorization boundary. **(D)** The similarity (Pearson correlation coefficient) between actual and template confusion matrices for each task and each ROI. Gray

dots represent individual participants, colored circles and error bars represent the mean \pm SEM across 7 participants. See Supplementary Figure 1 for an analogous analysis using a template for the *Nonlinear* task.

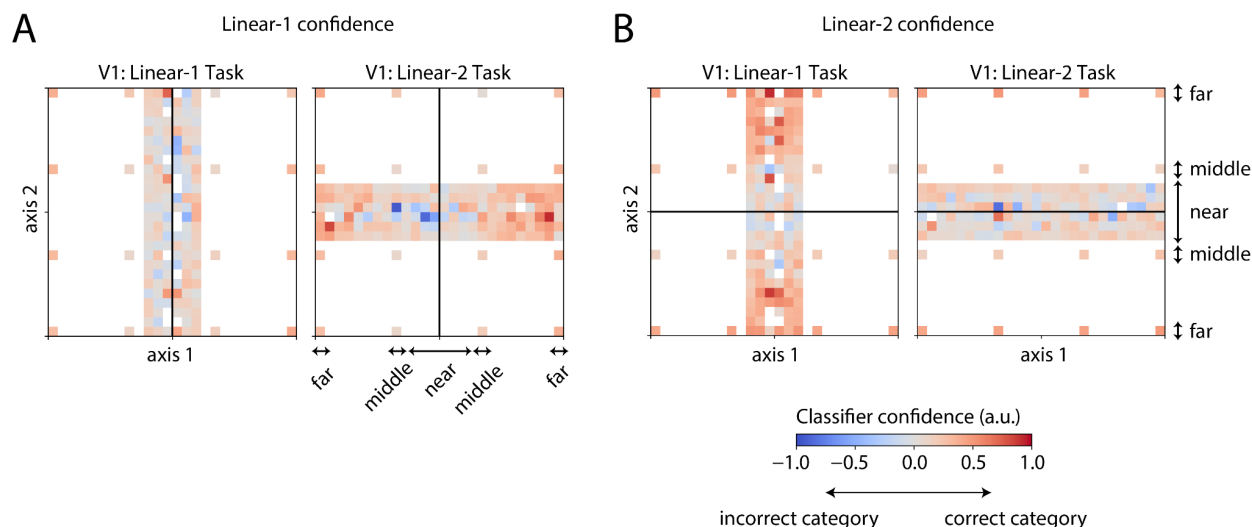


Figure 4. Illustration of how classifier “confidence” was computed with respect to each binary decision boundary. **(A)** *Linear-1* confidence, or confidence with respect to the *Linear-1* category boundary, was computed based on the difference between the total probability assigned by the 16-way classifier to each side of the boundary (see *Methods*). Left and right panels represent data from V1 in the *Linear-1* and *Linear-2* tasks, respectively, averaged across all participants. In each of the plots, each square represents a bin of shape space positions in the test dataset, and the color indicates the average confidence assigned to the correct category for that test trial (red) versus the incorrect category (blue). Arrows labeled “far”, “middle”, and “near” indicate the bins of distance from the category boundary in which confidence values were averaged (see Figure 5). **(B)** Same as A, but showing *Linear-2* confidence. An analogous procedure was also used to compute *Nonlinear* confidence; see *Methods*.

Next, we examined whether these representational changes differed depending on the position of shapes in the overall shape space. Specifically, we asked if changes in representations were more pronounced for shapes nearer to the category boundary than shapes further from the boundary. We divided the trials in each categorization task (*Linear-1*, *Linear-2*, and *Nonlinear* tasks) into three bins as a function of distance from each boundary: far, middle, and near (see Figure 4 and *Methods* for details). To measure the category separability of shapes in each of these distance bins, we computed classifier confidence with respect to each of the category boundaries: we refer to these measures as *Linear-1* confidence, *Linear-2* confidence, and *Nonlinear* confidence. Each type of confidence was computed by taking the output of the 16-way classifier described above and comparing the total probability assigned by the classifier to points on each side of each boundary. This analysis provides a continuous metric where larger positive values indicate higher separability of shapes across the boundary of interest.

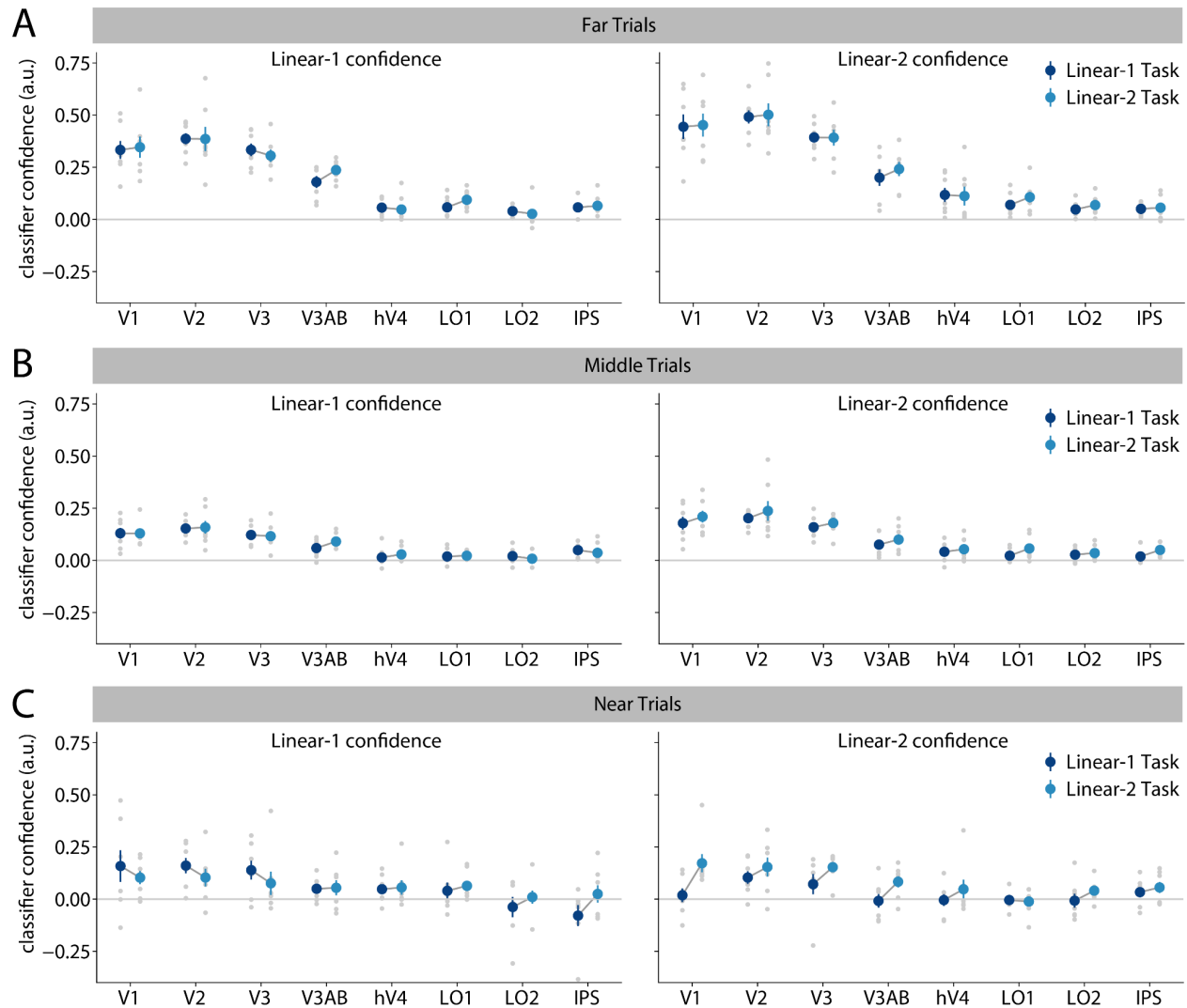


Figure 5. Discriminability of *Linear-1* and *Linear-2* shape categories depends on task and proximity to category boundaries. To obtain a continuous estimate of shape category discriminability, we used our 16-way multinomial classifier (see Figure 2) to compute classifier confidence toward the correct binary category on each trial (see Figure 4). Confidence was computed with respect to the *Linear-1* categorization boundary (*Linear-1* confidence; left) or the *Linear-2* categorization boundary (*Linear-2* confidence; right). **(A)** Confidence computed using “far” trials, meaning the 8 points in the main grid that fell furthest from the category boundary of interest. **(B)** Confidence computed using “middle” trials, meaning the 8 points in the main grid that fell nearest to the boundary of interest. **(C)** Confidence computed using “near” trials, meaning those that were not part of the main grid and fell nearest to the boundary of interest. The exact set of shape space positions sampled on near trials differed between tasks, but was matched using resampling for this analysis; see *Methods* for details. In **(A-C)**, the gray dots represent individual participants, colored circles and error bars represent the mean \pm SEM across 7 participants.

We first compared *Linear-1* confidence and *Linear-2* confidence across the *Linear-1* and *Linear-2* tasks. Overall, both types of confidence were highest for trials furthest from the boundary (Figure 5A), followed by middle trials (Figure 5B) and near trials (Figure 5C). This pattern is expected given that shapes further from the boundary are more distinctive from one another, while shapes nearer to the boundary are more ambiguous. In addition, this analysis revealed effects of task condition that differed for far, middle, and near trials. For trials in the far group, a three-way repeated measures ANOVA showed main effects of ROI and confidence boundary (i.e., *Linear-1* confidence versus *Linear-2* confidence), but no main effect of task or interaction between task and boundary (Supplementary Table 1), suggesting that discriminability of shapes across the *Linear-1* and *Linear-2* boundaries did not differ across tasks for this group of trials. For both the middle trials and the near trials, however, there was also a significant interaction between task and boundary (Supplementary Table 1). When each boundary was examined separately for each of these trial groups, we found a main effect of task on *Linear-2* confidence for both the middle trials and the near trials (two-way repeated measures ANOVA on middle trials; ROI: $F_{(7,42)} = 24.93$, $p < 0.001$; Task: $F_{(1,6)} = 10.22$, $p = 0.015$; ROI x Task: $F_{(7,42)} = 0.17$, $p = 0.990$; near trials; ROI: $F_{(7,42)} = 4.72$, $p < 0.001$; Task: $F_{(1,6)} = 10.25$, $p = 0.014$; ROI x Task: $F_{(7,42)} = 1.22$, $p = 0.315$), with *Linear-2* confidence showing higher values for the *Linear-2* task than the *Linear-1* task. Additional post-hoc tests in each ROI showed that for near trials, *Linear-2* confidence was significantly higher for the *Linear-2* task than the *Linear-1* task in V1 ($t_{(6)} = -3.72$, $p = 0.030$; paired t-test with permutation; see *Methods*). As with the confusion matrix analysis, the effect of task was larger for the *Linear-2* boundary than for the *Linear-1* boundary – there was no main effect of task seen for the *Linear-1* confidence values for either middle or near trials (middle trials; ROI: $F_{(7,42)} = 18.12$, $p < 0.001$, Task: $F_{(1,6)} = 0.10$, $p = 0.755$; ROI x Task: $F_{(7,42)} = 0.46$, $p = 0.868$; near trials; ROI: $F_{(7,42)} = 3.08$, $p = 0.006$; Task: $F_{(1,6)} = 0.01$, $p = 0.923$; ROI x Task: $F_{(7,42)} = 1.11$, $p = 0.379$).

In addition to comparing confidence across the two linear boundaries, we measured *Nonlinear* confidence for the far, middle, and near trials in each task (Figure 6). As before, confidence values tracked the distance of shapes from the boundary, with highest overall confidence observed for far trials, followed by middle and near trials. In contrast to the results with *Linear-2* confidence, however, *Nonlinear* confidence did not show any significant differences across tasks (Supplementary Table 2).

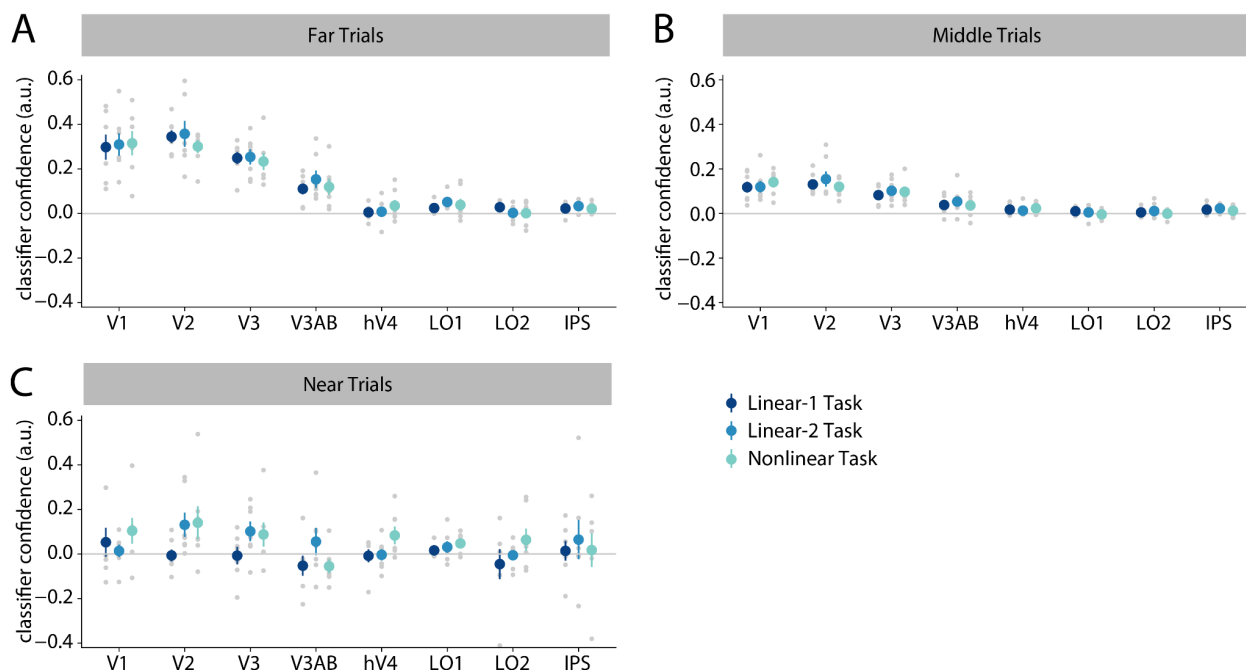


Figure 6. Discriminability of shapes across the *Nonlinear* boundary does not differ significantly across tasks. As in Figure 5, we computed the confidence of the classifier toward the correct *Nonlinear* task category for each trial. **(A)** Confidence computed using “far” trials, meaning the four points in the main grid that fell furthest from the two category boundaries (i.e., four corners of the shape space grid). **(B)** Confidence computed using “middle” trials, meaning the 12 points in the main grid that fell nearest to either of the two category boundaries. **(C)** Confidence computed using “near” trials, meaning those that were not part of the main grid and fell nearest to the two boundaries. The exact set of shape space positions sampled on near trials differed between tasks, but was matched using resampling for this analysis; see *Methods* for details. In **(A-C)**, the gray dots represent individual participants, colored circles and error bars represent the mean \pm SEM across 7 participants.

Finally, we evaluated whether the discriminability of shape representations across the relevant category boundary in each task was associated with behavioral performance. To test this, we compared classifier confidence for correct versus incorrect trials (focusing on “near” trials only, since these had the highest rate of incorrect responses). To ensure a fair comparison across correct and incorrect trials, we used bootstrap resampling to match the distribution of stimulus positions sampled in each group of trials; see *Methods* for details. As shown in Figure 7, this analysis revealed a significant difference in classifier confidence between correct and incorrect trials in both the *Linear-2* and the *Nonlinear* tasks, with confidence tending to be higher for correct trials than incorrect trials, particularly in early areas V1, V2, and V3. A two-way repeated measures ANOVA with factors of ROI and correctness revealed a significant main effect of correctness for both the *Linear-2* and *Nonlinear* tasks, and a significant interaction

between ROI x correctness for the *Nonlinear* task (*Linear-2*; ROI: $F_{(7,42)} = 5.57$, $p < 0.001$; Correctness: $F_{(1,6)} = 8.04$, $p = 0.0297$; ROI x Correctness $F_{(7,42)} = 1.96$, $p = 0.083$; *Nonlinear*; ROI: $F_{(7,42)} = 3.58$, $p = 0.004$; Correctness: $F_{(1,6)} = 8.05$, $p = 0.030$; ROI x Correctness $F_{(7,42)} = 3.41$, $p = 0.006$; p-values obtained using permutation test; see *Methods*). At the individual ROI level, confidence was significantly higher for correct than incorrect trials in V1 during the *Nonlinear* task ($t_{(6)} = 5.172$, $p = 0.018$; paired t-test with permutation; see *Methods*). The *Linear-1* task showed no significant differences in confidence for correct versus incorrect trials (ROI: $F_{(7,42)} = 4.45$, $p < 0.001$; Correctness: $F_{(1,6)} = 0.20$, $p = 0.664$; ROI x Correctness $F_{(7,42)} = 1.10$, $p = 0.382$). These results indicate that the separability of shape representations in early visual cortex across the task-relevant category boundary was associated with behavioral performance, at least for two out of three categorization tasks.

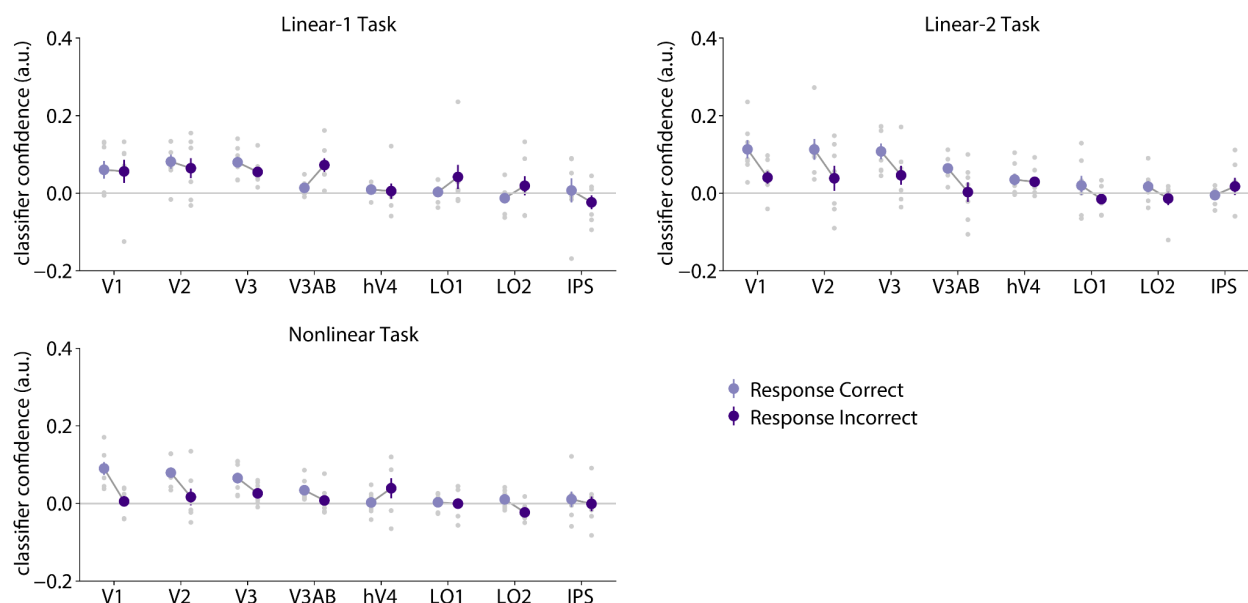


Figure 7: Task-relevant shape categories are more discriminable on correct versus incorrect trials. In each task, classifier confidence was computed with respect to the relevant category boundary for that task. Confidence was computed using “near” trials only (those nearest the relevant boundary), separately for trials with correct and incorrect behavioral responses. The set of shape space positions sampled on correct and incorrect trials was matched using resampling to ensure that the effect was not driven by stimulus differences; see *Methods* for details. Gray dots represent individual participants, colored circles and error bars represent the mean \pm SEM across 7 participants.

Discussion

Our goal was to determine whether and how human visual cortex representations of shape stimuli are adaptively modulated when switching between distinct task contexts. To test this, we trained participants to perform a categorization task on shape silhouette stimuli within a two-dimensional shape space (Figure 1). Participants categorized shapes according to different categorization rules (*Linear-1*, *Linear-2*, *Nonlinear*) on interleaved fMRI scanning runs, and we used multivariate decoding to explore how neural representations shift based on decision rules and the relative positions of shapes within the two-dimensional stimulus space. First, we used a confusion matrix analysis to show that shape representations became more aligned with the *Linear-2* boundary when participants were performing the *Linear-2* task versus the *Linear-1* task, with the largest effect observed in LO1 (Figure 3). We then showed that the discriminability of shapes across each linear boundary, as measured by classifier confidence, was higher when that boundary was relevant to the current task. These effects were most pronounced in early areas V1-V3, and were strongest for shapes located nearest to the active categorization boundary (Figure 5). Finally, we showed that the discriminability of shapes across relevant category boundaries was higher on correct versus incorrect trials, indicating a link with behavioral task performance (Figure 7). Together, these results demonstrate that performance of a categorization task with a dynamically changing task boundary is accompanied by changes to neural representations in human visual cortex.

The average accuracy of our decoder, across tasks, was highest in V2 followed by V1 and V3. This high decoding accuracy in early areas is surprising in light of earlier work suggesting that higher visual areas like ITC and LOC encode shapes similar to ours (i.e., radial frequency components (RFC)-defined silhouettes) in a way that matches perceptual similarity (Drucker & Aguirre, 2009; Op de Beeck et al., 2001), and that LOC is critically involved in shape computations (Vinberg & Grill-Spector, 2008). Work in non-human primates also indicates that neurons in ITC, as well as in V4, are more strongly tuned for shape and contour than neurons in V1 (Connor et al., 2007; DiCarlo & Maunsell, 2000; Pasupathy & Connor, 1999; Tanaka, 1993, 1996). One reason for our observation of higher decoding accuracy in early areas is that our stimuli were silhouettes presented at a fixed size and position, so invariance to size or position was not required to encode them accurately. As a result, fine-grained retinotopic and orientation tuning in areas like V1-V3 was likely sufficient to encode the shapes with high accuracy, without the need for an explicit – or invariant – contour or shape representation. Importantly, the goal of our experiment was not to measure abstract representations of shape or contour *per se* but to measure how visual representations change in accordance with dynamically varying decision boundaries, and our relatively simple stimulus set was appropriate for this goal.

The effects of task context on classifier confidence (Figure 5), as well as association of classifier confidence with behavioral performance (Figure 7), also tended to be strongest in early visual areas. This advantage for early areas may be due in part to the higher signal-to-noise ratio (SNR) of decoding accuracy in V1-V3, but it may also suggest that representations in these areas are particularly important for performance of our decision task. The findings of strong task-dependent effects in early retinotopic areas align with recent rodent studies, which

show that representations within sensory areas contain information pertinent to task goals, motor outcomes, and prior knowledge about sensory environments (Ebrahimi et al., 2022; Findling et al., 2023; Mimica et al., 2023; Niell & Stryker, 2010; Stringer et al., 2019). Extending these findings, our study demonstrates that human visual areas are more actively involved with decision-related computation than previously thought. Our results demonstrate that human sensory areas not only code for temporally varying task contexts but also dynamically integrate this information with incoming sensory inputs to optimize decision processes. This observation challenges the traditional view that sensory areas are primarily dedicated to basic sensory processing, suggesting a more multifaceted role in cognitive computation.

A plausible mechanism for guiding dynamic task coding and context-dependent representation of sensory inputs in humans may involve the deployment of selective attention. By flexibly prioritizing processing of relevant stimulus features based on current task goals, attention may guide the integration of sensory information with shifting task demands. Specifically, our observed task-dependent effects in early retinotopic areas are consistent with the literature on feature-based attention, which has shown that directing attention to simple visual features can modulate representations in early visual cortex (X. Chen et al., 2012; Foster & Ling, 2022; Gundlach et al., 2023; Jehee et al., 2011; Liu et al., 2003, 2007; Martinez-Trujillo & Treue, 2004; Mirabella et al., 2007; Saenz & Boynton, 2003; Serences & Boynton, 2007; Treue & Maunsell, 1996, 1999; Yoo et al., 2022). By modulating neurons coding for perceptual features that differentiate between categories, feature-based attention could provide a mechanism for improving the separability of different stimulus categories (Navalpakkam & Itti, 2007; Scolarì et al., 2012; Scolarì & Serences, 2009). Our result of early modulations is also consistent with Ester et al. (2020), who found biases in orientation representations that were related to categorization, although their paradigm used a single category boundary as opposed to a dynamically updated boundary. Additionally, other work using more complex stimuli such as three dimensional objects and human bodies has shown feature-based attention effects in higher visual areas such as LOC and the extrastriate body area (EBA), as opposed to early visual cortex (Jackson et al., 2017; Thorat & Peelen, 2022). These observations may indicate that attentional modulations in V1-V3 are more important for task performance when stimuli are relatively simple and require fine-grained spatial detail (e.g., oriented gratings, two-dimensional silhouettes), than when stimuli are more complex. In keeping with this idea of attention adapting dynamically to the most informative features for a task, a recent behavioral study demonstrated that feature-based attention is adaptively allocated according to experience with the variance of feature distributions (Witkowski & Geng, 2022). Our findings extend these prior studies by demonstrating feature-based attention as a potential mechanism for effectively integrating sensory information with changing task requirements within human sensory cortex.

Despite the relatively low classifier accuracy values that were observed in higher areas, we did observe a significant effect of task-relevance in LO1 based on the confusion matrix analysis in Figure 3. In this analysis, we demonstrated that classifier confusion matrices from LO1 were more aligned with the *Linear-2* task template during the *Linear-2* task versus the *Linear-1* task. The divergence of this finding from our classifier confidence analyses, in which early areas showed larger task effects than LO1, may indicate that the nature of representational changes in LO1 across categorization tasks differs from the changes in V1-V3.

Specifically, the confusion matrix analysis captures changes to the relationship between all 16 shapes in the main shape space grid, including pairs on the same side of the boundary, while the classifier confidence analysis only captures the discriminability of shapes across a single category boundary. One interesting possibility is that the changes in LO1 from the *Linear-1* task to the *Linear-2* task are primarily driven by re-structuring of shape representations within a given category (i.e. “acquired equivalence”; Goldstone, 1994) as opposed to an increase in discriminability across the boundary. Another possibility is that context-related changes in early areas reflect subtle changes in discriminability that allow the overall structure of the representational space to be largely maintained across tasks, while changes in LO1 reflect a more dramatic restructuring of sensory codes into a format that resembles a binary or categorical code for each task. Such a difference would be consistent with LO1 being a higher visual area more closely aligned with decision processes than early areas. Further experiments will be needed to evaluate these possibilities.

When classifier confidence values were broken down based on proximity to the category boundary, we observed the largest effects of categorization task on confidence for stimuli nearest the boundary, and no effect of task for the furthest stimulus positions. This scaling of categorization effects with proximity to the boundary is consistent with a previous fMRI experiment (Ester et al., 2020) as well as past behavioral experiments (Ashby & Maddox, 2005; Goldstone, 1994, 1998; Livingston et al., 1998; Newell & Bühlhoff, 2002). These convergent findings suggest that top-down modulatory effects in early visual cortex are strengthened on trials with higher category ambiguity, facilitating perceptual discrimination of these challenging stimuli. Importantly, our results also build on these past findings by demonstrating an increase in the discriminability of representations near the decision boundary during a task that requires flexible switching between multiple decision boundaries.

Task context had more consistent effects on discriminability with respect to the *Linear* tasks compared to the *Nonlinear* task, with no significant difference across tasks observed for *Nonlinear* confidence (Figure 6). This difference may be due to the fact that the *Nonlinear* task required using a non-linear decision boundary. The non-linear boundary was slightly more challenging behaviorally, as demonstrated by the slower RTs observed in the *Nonlinear* task compared to the *Linear-1* task, which is also consistent with a past report showing that a quadrant task with similar stimuli was more challenging for macaques to learn than a linear rule (Op de Beeck et al., 2001). The non-linearity of the boundary may also explain the lack of a consistent task-related modulation of *Nonlinear* confidence in visual cortex. It is possible that while top-down mechanisms are capable of selectively enhancing representations along one continuous axis in a perceptual space, such a mechanism does not exist for non-linear boundaries. On the other hand, it is also possible that practical aspects of the analysis prevented the effect from being measured. When *Nonlinear* confidence was computed for each task on “near” trials, we used only trials in the centermost region of the shape space grid (this is needed to facilitate a fair comparison across tasks; see *Methods*) and this resulted in a relatively small number of trials available for resampling. Given that the effect of task-relevance was expected to be strongest on near trials, such an effect might have been detectable if more trials were available. In support of this idea, though we did not observe a task-related modulation of *Nonlinear* confidence, we observed a significant within-task association of *Nonlinear* confidence

with behavioral performance (Figure 7). In this within-task analysis, a larger number of trials are available for resampling, which may have led to a more stable effect. Finally, the difference in outcomes between these analyses may also indicate that while discriminability of shapes across the *Nonlinear* boundary does not differ across task contexts, there is variability in the quality of representations across trials within the *Nonlinear* task, and this variability is associated with behavioral performance.

Comparing the two *Linear* tasks, we observed higher SNR for discriminating stimuli across the *Linear-2* boundary than the *Linear-1* boundary (i.e., higher values of similarity to *Linear-2* template and higher values of *Linear-2* confidence, across all tasks). We also observed more consistent effects of task relevance on *Linear-2* template similarity and *Linear-2* confidence than the analogous measures with respect to *Linear-1*. Finally, we did not observe any association of *Linear-1* confidence with behavioral performance, though such an effect was observed for *Linear-2* and *Nonlinear* confidence. These findings may be related to the difference in perceptual separability, as measured by our image similarity analyses, between the *Linear-1* and *Linear-2* categories. The *Linear-2* boundary, across which shapes are more perceptually distinctive, may also be a more effective target of context-dependent processing via selective attention mechanisms. Taken together, these findings may indicate an asymmetry in the allocation of attention to different dimensions within our shape space, in a way that reflects physical properties of the stimuli.

Overall, our findings provide evidence for context-dependent modulations of neural representations in early visual cortex, and show that these effects differ in accordance with temporally shifting task demands. Shape representations were modified to support discrimination of currently-relevant shape categories, with effects that were strongest for stimuli near the decision boundary. Moreover, these effects were associated with task performance. These results may indicate that visual cortex plays an active computational role in the flexible categorization of stimuli, providing new insight into how we organize knowledge about visual stimuli in the face of changing behavioral requirements.

Materials & Methods

Human participants

Seven (7) participants were recruited from the UCSD community, and were adults having normal or corrected-to-normal vision. Participants were between the ages of 24 and 32 (mean = 27.7, std = 3.0), and 4 out of 7 were female. The protocol for this study was approved by the Institutional Review Board at UCSD, and all participants provided written informed consent. As part of this experiment, each participant took part in one behavioral training session lasting approximately 1 hour, for which they were compensated at a rate of \$10/hour and three scanning sessions each lasting approximately 2 hours, for which they were compensated at a rate of \$20/hour. During each scanning session for this experiment, participants also performed several runs of a n-back (repeat detection) task on the same stimuli used in our main task (see *Main task design*). Data from this task are not analyzed here but are included in our full open dataset (see *Data availability*). Each participant also participated in a separate retinotopic mapping scan session; for five participants this retinotopic mapping session was performed as part of an earlier experiment and for the remaining two it was performed just prior to the start of the present experiment.

Acquisition of MRI data

All magnetic resonance imaging (MRI) scanning was performed using a General Electric (GE) Discovery MR750 3.0T research-dedicated scanner at the UC San Diego Keck Center for Functional Magnetic Resonance Imaging. We used a Nova Medical 32-channel head coil (NMSC075-32-3GE-MR750) to acquire all functional echo-planar imaging (EPI) data, using the Stanford Simultaneous Multislice (SMS) EPI sequence (MUX EPI), with a multiband factor of 8 and 9 axial slices per band (total slices = 72; 2 mm³ isotropic; 0 mm gap; matrix = 104 x 104; field of view = 20.8 cm; repetition time/time to echo [TR/TE] = 800/35 ms; flip angle = 52°; inplane acceleration = 1). To perform image reconstruction and un-aliasing we used reconstruction code from the Stanford Center for Neural Imaging, on servers hosted by Amazon Web Services. The initial 16 TRs collected at sequence onset were used as reference images in order to transform data from k-space to image space. In addition, two “topup” datasets (17s each) were collected at the halfway point of each session, using forward and reverse phase-encoding directions. These runs were used to correct for distortions in the EPI sequences from the same session using topup functionality (Andersson et al., 2003) in the FMRIB Software Library (FSL; Jenkinson et al., 2012).

In addition to the functional data, we also collected a high-resolution anatomical scan for each participant as part of that participant’s retinotopic mapping session. This anatomical T1 image was used for segmentation, flattening, and delineation of the retinotopic mapping data. For three out of the seven participants, we acquired this anatomical scan using the same 32 channel head coil used for functional scanning, and for the remaining four participants, we used

an in vivo eight-channel head coil. Anatomical scans were acquired using accelerated parallel imaging (GE ASSET on a FSPGR T1-weighted sequence; $1 \times 1 \times 1 \text{ mm}^3$; 8136 ms TR; 3172 ms TE; 8° flip angle; 172 slices; 1 mm slice gap; 256 x 192 cm matrix size). When the 32-channel head coil was used, anatomical scans were corrected for inhomogeneities in signal intensity using GE's 'phased array uniformity enhancement' (PURE) method.

Preprocessing of functional MRI data

Preprocessing of functional data was performed using tools from FSL and FreeSurfer (available at <http://www.fmrib.ox.ac.uk/fsl> and <https://surfer.nmr.mgh.harvard.edu>). We first performed cortical surface gray-white matter volumetric segmentation of the anatomical T1 scans for each participant, using the recon-all function in FreeSurfer (Dale et al., 1999). The segmented T1 data were then used to define cortical meshes on which we defined retinotopic ROIs (see next section for details). We also used the anatomical T1 data in order to align multi-session functional data to a common space for each participant. This was performed by using the first volume of the first scan for each session as a template, and using this template to align the entire functional session to the anatomical scan for each participant. We used the manual and automatic boundary-based registration tools in FreeSurfer to perform co-registration between functional and anatomical data (Greve & Fischl, 2009), then used the resulting transformation matrix and FSL FLIRT to transform all functional data into a common space (Jenkinson et al., 2002; Jenkinson & Smith, 2001). Next, we used FSL MCFLIRT to perform motion correction (Jenkinson et al., 2002), with no spatial smoothing, with a final sinc interpolation stage, and 12° of freedom. Finally, we performed de-trending to remove slow drifts in the data using a high-pass filter (1/40 Hz cutoff).

Following these initial preprocessing stages, we z-scored the data within each scan run on a per-voxel basis to correct for differences in mean and variance across runs. This and all subsequent analyses were performed using Python 3.7.10 (Python Software Foundation, Wilmington, DE). Next, we obtained a single estimate for each voxel's activation on each trial by averaging the time series over a window spanning from 4-7 TRs following image onset (see *Main task design* for more details on task timing and procedure). We then extracted data from voxels within several regions of interest (ROIs; see next section) that were used for subsequent analyses.

Retinotopic ROI definitions

We defined several retinotopic visual ROIs: V1, V2, V3, V3AB, hV4, LO1, LO2, and IPS, following previously identified retinotopic mapping procedures (Engel et al., 1997; Jerde & Curtis, 2013; Sereno et al., 1995; Swisher et al., 2007; Wandell et al., 2007; Winawer & Witthoft, 2015). We combined all intraparietal sulcus (IPS) subregions (IPS0, IPS1, IPS2, IPS3), into a single combined IPS ROI, as this led to improved classifier accuracy relative to the individual sub-regions. During retinotopic mapping runs, participants viewed black-and-white contrast reversing checkerboard stimuli that were configured as a rotating wedge (10 cycles, 36 s/cycle), expanding ring (10 cycles, 32 s/cycle), or bowtie shape (8 cycles, 40 s/cycle). During the

rotating wedge task, a contrast detection task (detecting dimming events approximately every 7.5 s) was used to encourage covert attention to the stimulus. Average accuracy on this task was $79.07 \pm 3.84\%$ (mean \pm SEM across 7 participants). The stimulus had a maximum eccentricity of 9.3° .

After defining retinotopic ROIs using these methods, we further thresholded the ROIs using an independent localizer task to identify voxels that were responsive to the region of space in which shape stimuli could appear (see *Silhouette localizer task* for details on this task). The data from the localizer were analyzed using a general linear model (GLM) implemented in FSL's FMRI Expert Analysis Tool (FEAT; version 6.00). This analysis included performing brain extraction and pre-whitening (Smith, 2002; Woolrich et al., 2001). We generated predicted BOLD responses by convolving each stimulus onset with a canonical gamma hemodynamic response (phase = 0s, s.d. = 3s, lag = 6s), and combined individual runs using a standard weighted fixed effects analysis. We identified voxels that were significantly activated by the stimulus versus baseline ($p < 0.05$, false discovery rate (FDR) corrected). This mask of responsive voxels was then intersected with each ROI definition to obtain the final thresholded ROI definitions. The exception to this was the IPS ROIs, to which we did not apply any additional thresholding; this was because the localizer yielded few responsive voxels in IPS for some participants.

Shape stimuli

We used a set of shape silhouette stimuli that varied parametrically along two continuous dimensions, generating a 2-dimensional shape space (Figure 1A). Each shape in this space was a closed contour composed of radial frequency components (RFCs; Op de Beeck et al., 2001; Zahn & Roskies, 1972). Each shape was composed of 7 different RFCs, where each component has a frequency, amplitude, and phase. To generate the 2-dimensional shape space, we parametrically varied the amplitude of two RFCs, leaving the others constant. The manipulation of RFC amplitude was used to define an x/y grid in arbitrary units that spanned positions between 0-5 a.u., with adjacent grid positions spaced by 0.1 a.u. All shape space positions on all trials were sampled from this grid of shape space positions. We also defined a coarser grid of 16 points (a 4x4 grid) which was used to generate the 16 stimuli that were shown on the majority of trials; this grid is referred to as the "main grid", and included all x/y combinations of the points [0.1, 1.7, 3.3, 4.9] in shape space coordinates. Stimuli corresponding to points in shape space that were not part of the main grid were used to make the tasks more difficult, see *Main task design* for details.

We divided the shape space into four quadrants by imposing boundaries at the center position of the grid (2.5 a.u.) in each dimension. To define the binary categories that were relevant for each task (see *Main task design*), we grouped together two quadrants at a time, with the *Linear-1* task and *Linear-2* tasks grouping quadrants that were adjacent (creating either a vertical or horizontal linear boundary in shape space), and the *Nonlinear* task grouping quadrants that were non-adjacent (creating a non-linear boundary). During task training as well as before each scanning run, we utilized a "prototype" image for each shape space quadrant as

a way of reminding participants of the current categorization rule. The prototype for each quadrant was positioned directly in the middle of the four main grid positions corresponding to that quadrant (i.e. the x/y coordinates for the prototypes were combinations of [0.9, 4.1] a.u.). These prototype images were never shown during the categorization task trials, to prevent participants from simply memorizing the prototypes. Shapes used in the task were also never positioned exactly on any quadrant boundary in order to prevent any ambiguity about category.

Display parameters

During all scanning runs, stimuli were presented to participants by projecting onto a screen that was mounted on the inside of the scanner bore, just above the participant's chest. The screen was visible to the participant via a mirror that was attached to the head coil. The image projected onto the screen was a rectangle with maximum horizontal eccentricity of 13 degrees (center-to-edge distance) and maximum vertical eccentricity of 10 degrees. In the main task and silhouette localizer task, the region of the screen in which shapes could appear subtended a maximum eccentricity of 11 degrees in the horizontal direction, and 9 degrees in the vertical direction. The fixation point in all tasks was a gray square 0.2 degrees in diameter; participants were instructed to maintain fixation on this point throughout all experimental runs.

In the main task, shapes were displayed as gray silhouettes on a gray background. For all participants except for the first participant (S01), the shapes were darker than the background (shape = 31, background = 50; luminance values are in the range 0-255). For S01, the shapes were lighter than the background (shape = 230, background = 77). The change in parameters was made because the brighter stimuli shown to S01 led to display artifacts when scanning subsequent participants, and darker stimuli reduced these artifacts. S01 reported no artifacts and performed well on the task. No gamma correction was performed.

Main task design

The main experimental task consisted of categorizing shape silhouette stimuli (Figure 1) into binary categories. There were three task conditions: *Linear-1*, *Linear-2*, and *Nonlinear*, each of which corresponded to a different binary categorization rule. Shape stimuli were drawn from a two-dimensional shape space coordinate system (see *Shape stimuli*). The *Linear-1* and *Linear-2* tasks used a boundary that was linear in this shape space, while the *Nonlinear* task used a boundary that was non-linear in this shape space (requiring participants to group non-adjacent quadrants into a single category, see Figure 1 for illustration). Each trial consisted of the presentation of one shape for 1s, and trials were separated by an inter-trial interval (ITI) that was variable in length, uniformly sampled from the interval 1-5s. Participants responded on each trial with a button press (right index or middle finger) to indicate which binary category the currently viewed shape fell into; the mapping between category and response was counter-balanced within each scanning session. Participants were allowed to make a response anytime within the window of 2s from stimulus onset. Feedback was given at the end of each run, and

included the participant's overall accuracy, as well as their accuracy broken down into "easy" and "hard" trials (see next paragraph for description of hard trials), and the number of trials on which they failed to respond. No feedback was given after individual trials.

Each run in the task consisted of 48 trials and lasted 261s (327 TRs). Of the 48 trials, 32 of these used shapes that were sampled from a grid of 16 points evenly spaced within shape space ("main grid", see *Shape stimuli*), each repeated twice. These 16 shapes were presented twice per run regardless of task condition. The remaining 16 trials (referred to as "hard" trials) used shapes that were variable depending on the current task condition and the difficulty level set by the experimenter. The purpose of these trials was to allow the difficulty level to be controlled by the experimenter so that task accuracy could be equalized across all task conditions, and prevent any single task from being trivially easy for each participant. For each run of each task, the experimenter selected a difficulty level between 1-13, with each level corresponding to a particular bin of distances from the active categorization boundary (higher difficulty denotes closer distance to boundary). For the *Nonlinear* task, the distance was computed as a linear distance to the nearest boundary. The "hard" trials were generated by randomly sampling 16 shapes from the specified distance bin, with the constraint that 4 of the shapes had to come from each of the four quadrants in shape space. This manipulation ensured that responses were balanced across categories within each run. For many of the analyses presented here, we excluded these hard trials, focusing only on the "main grid" trials where the same images were shown across all task conditions.

Participants performed 12 runs of the main task within each scanning session, for a total of 36 runs across all 3 sessions (with the exception of one participant (S06) for whom 3 runs are missing due to a technical error). The 12 runs in each session were divided into 6 total "parts" where each part consisted of a pair of 2 runs having the same task condition and the same response mapping (3 conditions x 2 response mappings = 6 parts). Each part was preceded by a short training run, which consisted of 5 trials, each trial consisting of a shape drawn from the main grid. The scanner was not on during these training runs, and the purpose of these was to remind the participant of both the currently active task and the response mapping before they began performing the task runs for that part. The order in which the 6 parts were shown was counter-balanced across sessions. Before each scan run began, the participant was again reminded of the current task and response mapping via a display that presented four prototype shapes, one for each shape space quadrant (see *Shape stimuli* for details on prototype shapes). The prototypes were arranged with two to the left of fixation and two to the right of fixation, and the participant was instructed that the two leftmost shapes corresponded to the index finger button and the two rightmost shapes corresponded to the middle finger button. This display of prototype shapes was also used during the training runs to provide feedback after each trial: after each training trial, the four prototype shapes were shown, and the two prototypes corresponding to the correct category were outlined in green, with accompanying text that indicated whether the participant's response was correct or incorrect. This feedback display was not shown during the actual task runs.

Before the scan sessions began, participants were trained to perform the shape categorization tasks in a separate behavioral session (training session took place on average 4.7 days before the first scan session). During this behavioral training session, participants performed the same task that they performed in the scanner, including 12 main task runs (2 runs for each combination of condition and response mapping; i.e., each of the 6 parts). As in the scan sessions, each part was preceded by training runs that consisted of 5 trials, each accompanied by feedback. Participants completed between 1-3 training runs before starting each part. Average training session accuracy was 0.82 ± 0.01 (mean \pm SEM across participants) for the *Linear-1* task, 0.81 ± 0.02 for the *Linear-2* task, and 0.78 ± 0.03 for the *Nonlinear* task.

Silhouette localizer task

A silhouette localizer task was used to identify voxels that were responsive to all the regions of retinotopic space where the shape stimuli could appear. For this task, a single silhouette shape was generated that covered the area spanned by any shape in the main grid. The silhouette region was rendered with a black-and-white flashing checkerboard (spatial period = 2 degrees) against a mid-gray background. On each trial, the flashing checkerboard silhouette stimulus appeared for a total duration of 7s, with trials separated by an ITI that varied between 2-8s (uniformly sampled). During each trial the checkerboard was flashed with a frequency of 5 Hz (1 cycle = on for 100 ms, off for 100 ms). On each cycle, the checkerboard was re-drawn with a randomized phase. There were 20 trials per run of this task, and participants performed between 4 and 7 runs of this task across all sessions. During all runs of this task, participants were instructed to monitor for a contrast dimming event and press a button when the dimming occurred. Dimming events occurred with a probability of 0.10 on each frame, and were separated by a minimum of 4 cycles. There were on average 17 dimming events in each run (minimum 10; maximum 25). Average hit rate (proportion of events correctly detected) was 0.72 ± 0.10 (mean \pm s.d. across participants), and the average number of false alarms per run was 1.98 ± 2.05 (mean \pm s.d. across participants).

Image similarity analysis

To estimate the perceptual discriminability of our shape categories, we used two computer vision models to extract activations in response to each stimulus image. We first used the GIST model (Oliva & Torralba, 2001), which is based on Gabor filters and captures low-level spectral image properties. We also extracted features from a pre-trained SimCLR model (T. Chen et al., 2020), which is a self-supervised model trained using contrastive learning on a large image database. We selected these two models because the GIST model captures clearly defined image properties similar to those represented in the early visual system, while the SimCLR model can capture a wider set of image features, including mid-level and high-level properties. The GIST model was implemented in Matlab, using a 4x4 spatial grid, 4 spatial scales, and 4 orientations per spatial scale. The version of SimCLR that we used was

implemented in PyTorch and used a ResNet-50 backbone (pre-trained model downloaded from <https://pypi.org/project/simclr/>). We extracted activations from blocks [2,6,12,15] and performed a max-pooling operation (kernel size = 4, stride = 4) to reduce the size of activations from each block. We used principal components analysis (PCA) to further reduce the size of activations, retaining a maximum of 500 components per block, and concatenated the resulting features across all blocks.

Using these activations, we computed the separability of shape categories across each of our boundaries (*Linear-1*, *Linear-2*, *Nonlinear*) by computing all pairwise Euclidean distances between main grid shapes in the same category (within-category distances) and main grid shapes in different categories (between-category distances). We then computed the average of the within-category distances (w) and between-category distances (b). The separability measure for each boundary was computed as: $(b-w)/(b+w)$.

Multivariate classifier analysis

We used a multivariate classifier to estimate how well the voxel activation patterns from each ROI could be used to discriminate different shape stimuli. Classification was performed within each participant, each ROI, and each task condition separately. Before training the classifier, we mean-centered the activation patterns on each trial, by subtracting the average signal across voxels from each trial. We cross-validated the classifier by leaving one run out at a time during training, looping over held-out test runs so that every run served as the test run once. During training of the classifier, we used only trials on which main grid shapes were shown, meaning there were 16 unique shapes that were treated as distinct classes. We then constructed a 16-way multinomial logistic regression classifier, implemented using *scikit-learn* (version 1.0.2) in Python 3.6. We used the 'lbfgs' solver and L2 regularization. To select the L2 regularization parameter (C), we created a grid of 20 candidate C values that were logarithmically spaced between 10^{-9} and 1. We then used nested cross-validation on the training data only to select the C resulting in highest accuracy across folds, and re-fit the model for the entire training set using the best C parameter. The resulting classifier was then used to predict the class (1-16) for all trials in the test dataset (note that this included some trials where the viewed shape was not in the main grid, and thus was not included in classifier training). In addition to a predicted class for each trial, the classifier returned a continuous probability estimate for each of the 16 classes, obtained using a softmax function.

To evaluate whether the accuracy of the classifier was significantly greater than chance, we used a permutation test. To do this, we performed 1000 iterations of training and testing the classifier, constructed in the same way as described above, using shuffled labels for the data. We always performed shuffling within a given scan run, so that the run labels were kept intact, and leave-run-out cross-validation was performed as in the original method. To make this computationally feasible, we did not perform C selection on every shuffling iteration, instead we used a fixed C value of 0.023, which was approximately the median of the C values obtained across all models fit to the real data. We obtained a p-value for each individual participant, ROI, and task condition by computing the proportion of shuffle iterations on which shuffled classifier

accuracy was greater than or equal to the real classifier accuracy. To obtain p-values for the participant-averaged classification accuracy for each ROI and task, we used the same procedure but first averaged the values across participants, within each shuffle iteration. All reported p-values were false-discovery-rate (FDR) corrected at $q = 0.01$ (Benjamini & Hochberg, 1995).

Confusion matrix analysis

For each participant, ROI, and task, we generated a confusion matrix for the 16-way multinomial classifier. This was a 16 x 16 matrix where each row represents the set of trials on which a given shape was actually shown, and each column in the row represents the proportion of those trials that the classifier assigned into each of the 16 classes, and each row sums to 1. To compute confusion matrices we used only trials in the main grid, and only used trials on which the participant made a correct behavioral response. To quantify the alignment of confusion matrices with the representation needed to solve each task, we generated template confusion matrices for the *Linear-1* and *Linear-2* tasks, where each template matrix had 0 for pairs of stimuli that were on different sides of the boundary and 1 for pairs of stimuli that were on the same side of the boundary. We then computed the Pearson correlation coefficient between each actual confusion matrix and each template confusion matrix.

Classifier confidence

For several analyses, we were specifically interested in the discriminability of shapes belonging to different binary categories. To measure the discriminability of shapes across each boundary (*Linear-1*, *Linear-2*, *Nonlinear*), we used the continuous probability estimates output by the 16-way classifier to compute classifier confidence with respect to each boundary. For each boundary and each individual trial, our measure of classifier confidence was computed as the difference between the total probability assigned by the classifier to the “correct” binary category for that trial [$p(\text{correct})$] and the total probability assigned by the classifier to the “incorrect” binary category for that trial [$p(\text{incorrect})$]. We obtained $p(\text{correct})$ by summing the probability assigned to the 8 main grid shapes in the same category as the shape on the current trial (based on whichever category boundary was currently being considered), and $p(\text{incorrect})$ by summing the probability assigned to the 8 main grid shapes in the other category. Note that this measure of confidence can be computed even when the test trial shape is not part of the main grid. To interpret this measure, large positive values of confidence indicate high discriminability of shapes across a given category boundary, and large negative or zero values indicate poor discriminability.

For the analyses where confidence values are broken down by “far”, “middle” and “near” trials, the far and middle trials are always restricted to positions in the main grid. For the *Linear-1* and *Linear-2* tasks, there are 8 total positions counted as far and 8 counted as middle. For the *Nonlinear* task, we counted the 4 corner positions as far and the 12 other positions as middle. The near trials are always points that are not part of the main grid; see next section for details on how these points were sampled to compute average confidence. When average confidence

values are reported, they are averaged over behaviorally correct trials only (unless otherwise specified).

Bootstrap resampling procedures

When comparing classifier confidence values across tasks on “near” trials (i.e. those closest to each boundary and not in the main grid), we used bootstrap resampling to match the distribution of shape positions sampled in each task for each participant. This was implemented because the range of shape positions that were sampled on near trials differed between tasks (see Figure 1B), and this difference in stimulus properties could have, if not corrected, contributed to a difference in average confidence across tasks. Resampling was performed with respect to one categorization boundary at a time and for each participant separately. For each of the linear boundaries (*Linear-1* and *Linear-2*), we used resampling to equate the set of positions sampled between the *Linear-1* and *Linear-2* tasks. To achieve this, for each boundary we collapsed the set of coordinates sampled on the near trials in each task onto a single axis that ran perpendicular to the boundary of interest. We then binned the coordinates into a set of 6 linearly-spaced bins that spanned the portion of shape space nearest the boundary (from 1.8 to 3.2 in shape space coordinates; see *Shape stimuli*). Since not all bins were necessarily sampled in both *Linear-1* and *Linear-2* tasks (this depended on the task difficulty level for each participant), we then selected a subset of bins that were 1) sampled from in both task conditions and 2) were also symmetric around the categorization boundary (this could be all 6 bins, 4 bins, or 2 bins). For each task, we then performed 1,000 iterations on which we resampled with replacement a set of approximately 50 trials that evenly sampled from each bin, and computed the average classifier confidence for this resampled set. The final confidence values for each participant reflect the average across these 1,000 bootstrapping iterations.

When computing confidence with respect to the *Nonlinear* boundary, we used the same procedure to equate the set of positions sampled between all three categorization tasks (*Linear-1*, *Linear-2*, *Nonlinear*). To bin positions in this case, instead of collapsing coordinates onto a single axis, we computed the distance between each [x,y] coordinate and the nearest linear boundary, and multiplied by (+1) for coordinates in nonlinear category 1 or (-1) for coordinates in nonlinear category 2, which results in a single coordinate value that captures distance from the boundary as well as category sign. We restricted the set of included points to those that fell within the centermost square region of the grid (from 1.8 to 3.2 in shape space coordinates along both axes). These coordinate values were then binned and resampled as in the original procedure.

We also used bootstrap resampling to equate the distribution of coordinate positions sampled on correct versus incorrect trials. This resampling was always performed within one task at a time, and only for confidence values with respect to the task-relevant categorization axis. As in the other procedure, we binned the sampled coordinates along the relevant axis into a set of linearly-spaced bins (between 1.8 and 3.2 in shape space coordinates). In this case we used 12 bins, because the sampling of points was more dense when considering the task-relevant axis only. We then identified a subset of these 12 bins that were sampled on both

correct and incorrect trials, and were also symmetric around the categorization boundary. We then resampled with replacement a set of approximately 100 correct trials and approximately 100 incorrect trials that each evenly sampled from all bins, averaged the confidence values across these 100 trials, and repeated this procedure 1000 times.

Statistical analysis

To perform statistical comparisons of classifier confidence values and template correlation coefficient values (see previous sections) across ROIs and categorization tasks, we used repeated measures ANOVA tests, implemented using *statsmodels* in Python 3.6. To obtain non-parametric p-values for these tests (which are suitable for small sample sizes), we performed permutation tests where we shuffled the values within each participant 10,000 times, and computed F-statistics for each effect on the shuffled data. This resulted in a null distribution of F-values for each effect. The final p-values for each effect were based on the proportion of iterations on which the shuffled F-statistic was greater than or equal to the real F-statistic. F-statistics reported in the text reflect those obtained using the real (unshuffled) data.

To perform post-hoc tests for differences between tasks in each ROI, we used a paired t-test with permutation. For each ROI, we computed a t-statistic for the true difference between the conditions across participants, then performed 10,000 iterations where we randomly swapped the values within each participant across conditions, with 50% probability. This resulted in a null distribution of t-statistics. The final two-tailed p-value was obtained by computing the proportion of iterations on which the shuffled t-statistic was greater than or equal to the real t-statistic and the proportion of iterations on which the real t-statistic was greater than or equal to the shuffled t-statistic, taking the minimum and multiplying by 2.

Code availability statement

All code required to reproduce our analyses is available at <https://github.com/mmhenderson/shapeDim>.

Data availability

All data used in the present study will be deposited as MATLAB-formatted data in Open Science Framework.

Acknowledgments

This work was supported by NEI R01-EY025872 to JS, NIMH Training Grant in Cognitive Neuroscience (T32-MH020002) to MH, the Swartz Foundation Fellowship for Theory in Neuroscience to NR, and the Kavli Institute for Brain and Mind Postdoctoral Award to NR. We thank Stephanie Nelli for helpful discussions during the inception of this project as well as Anna

Shafer-Skelton and Julie Eitzen for their help with data collection. We would like to acknowledge Hans Op de Beeck for providing code that was used for stimulus generation.

Author Contributions

MH and JS conceived the research. MH, JS, and NR designed, performed the research, analyzed data, and wrote the manuscript.

Declaration of Interests

The authors declare no competing interests.

References

- Andersson, J. L. R., Skare, S., & Ashburner, J. (2003). How to correct susceptibility distortions in spin-echo echo-planar images: Application to diffusion tensor imaging. *NeuroImage*, 20(2), 870–888. [https://doi.org/10.1016/S1053-8119\(03\)00336-7](https://doi.org/10.1016/S1053-8119(03)00336-7)
- Ashby, F. G., & Maddox, W. T. (2005). Human category learning. *Annual Review of Psychology*, 56, 149–178. <https://doi.org/10.1146/ANNUREV.PSYCH.56.091103.070217>
- Benjamini, Y., & Hochberg, Y. (1995). Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society: Series B (Methodological)*, 57(1), 289–300. <https://doi.org/10.1111/J.2517-6161.1995.TB02031.X>
- Biederman, I., Glass, A. L., & Stacy, E. W. (1973). Searching for objects in real-world scenes. *Journal of Experimental Psychology*, 97(1), 22–27. <https://doi.org/10.1037/h0033776>
- Bruner, J. S. (1957). On perceptual readiness. *Psychological Review*, 64(2), 123–152. <https://doi.org/10.1037/h0043805>
- Chen, T., Kornblith, S., Norouzi, M., & Hinton, G. (2020). *A Simple Framework for Contrastive Learning of Visual Representations*.
- Chen, X., Hoffmann, K.-P., Albright, T. D., & Thiele, A. (2012). Effect of feature-selective

- attention on neuronal responses in macaque area MT. *Journal of Neurophysiology*, 107(5), 1530–1543. <https://doi.org/10.1152/jn.01042.2010>
- Connor, C. E., Brincat, S. L., & Pasupathy, A. (2007). Transformation of shape information in the ventral pathway. *Current Opinion in Neurobiology*, 17(2), 140–147. <https://doi.org/10.1016/j.conb.2007.03.002>
- Cromer, J. A., Roy, J. E., & Miller, E. K. (2010). *Representation of Multiple, Independent Categories in the Primate Prefrontal Cortex*. <https://doi.org/10.1016/j.neuron.2010.05.005>
- Dale, A. M., Fischl, B., & Sereno, M. I. (1999). Cortical surface-based analysis: I. Segmentation and surface reconstruction. *NeuroImage*, 9(2), 179–194. <https://doi.org/10.1006/nimg.1998.0395>
- DiCarlo, J. J., & Maunsell, J. H. R. (2000). Form representation in monkey inferotemporal cortex is virtually unaltered by free viewing. *Nature Neuroscience*, 3(8), Article 8. <https://doi.org/10.1038/77722>
- Drucker, D. M., & Aguirre, G. K. (2009). Different spatial scales of shape similarity representation in lateral and ventral LOC. *Cerebral Cortex (New York, N.Y. : 1991)*, 19(10), 2269–2280. <https://doi.org/10.1093/cercor/bhn244>
- Ebrahimi, S., Lecoq, J., Rummyantsev, O., Tasci, T., Zhang, Y., Irimia, C., Li, J., Ganguli, S., & Schnitzer, M. J. (2022). Emergent reliability in sensory cortical coding and inter-area communication. *Nature* 2022 605:7911, 605(7911), 713–721. <https://doi.org/10.1038/s41586-022-04724-y>
- Engel, S., Glover, G. H., & Wandell, B. A. (1997). Retinotopic organization in human visual cortex and the spatial precision of functional MRI. *Cerebral Cortex*, 7(2), 181–192. <https://doi.org/10.1093/cercor/7.2.181>
- Ester, E. F., Sprague, T. C., & Serences, J. T. (2020). Categorical biases in human occipitoparietal cortex. *Journal of Neuroscience*, 40(4), 917–931.

<https://doi.org/10.1523/JNEUROSCI.2700-19.2019>

- Findling, C., Hubert, F., Laboratory, I. B., Acerbi, L., Benson, B., Benson, J., Birman, D., Bonacchi, N., Carandini, M., Catarino, J. A., Chapuis, G. A., Churchland, A. K., Dan, Y., Dewitt, E. E., Engel, T. A., Fabbri, M., Faulkner, M., Fiete, I. R., Freitas-Silva, L., ... Pouget, A. (2023). Brain-wide representations of prior information in mouse decision-making. *bioRxiv*, 2023.07.04.547684. <https://doi.org/10.1101/2023.07.04.547684>
- Folstein, J., Palmeri, T. J., Van Gulick, A. E., & Gauthier, I. (2015). Category Learning Stretches Neural Representations in Visual Cortex. *Current Directions in Psychological Science*, 24(1), 17–23. <https://doi.org/10.1177/0963721414550707>
- 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. <https://doi.org/10.1093/cercor/bhs067>
- Foster, J. J., & Ling, S. (2022). Feature-based attention multiplicatively scales the fMRI-BOLD contrast-response function. *The Journal of Neuroscience: The Official Journal of the Society for Neuroscience*, 42(36), 6894–6906. <https://doi.org/10.1523/JNEUROSCI.0513-22.2022>
- Freedman, D. J., & Assad, J. A. (2016). Neuronal Mechanisms of Visual Categorization: An Abstract View on Decision Making. *Annual Review of Neuroscience*, 39(1), 129–147. <https://doi.org/10.1146/annurev-neuro-071714-033919>
- Freedman, D. J., Riesenhuber, M., Poggio, T., & Miller, E. K. (2003). A comparison of primate prefrontal and inferior temporal cortices during visual categorization. *The Journal of Neuroscience: The Official Journal of the Society for Neuroscience*, 23(12), 5235–5246. <https://doi.org/10.1523/JNEUROSCI.2312-03.2003> [pii]
- Gauthier, I., Skudlarski, P., Gore, J. C., & Anderson, A. W. (2000). Expertise for cars and birds recruits brain areas involved in face recognition. *Nature Neuroscience*, 3(2), 191–197. <https://doi.org/10.1038/72140>

- Goldstone, R. (1994). Influences of Categorization on Perceptual Discrimination. *Journal of Experimental Psychology: General*, 123(2), 178–200. <https://doi.org/10.1037/0096-3445.123.2.178>
- Goldstone, R. (1998). Perceptual learning. *Annual Review of Psychology*, 49, 585–612. <https://doi.org/10.1146/ANNUREV.PSYCH.49.1.585>
- Greve, D. N., & Fischl, B. (2009). Accurate and robust brain image alignment using boundary-based registration. *NeuroImage*, 48(1), 63–72. <https://doi.org/10.1016/j.neuroimage.2009.06.060>
- Gundlach, C., Wehle, S., & Müller, M. M. (2023). Early sensory gain control is dominated by obligatory and global feature-based attention in top-down shifts of combined spatial and feature-based attention. *Cerebral Cortex (New York, N.Y.: 1991)*, bhad282. <https://doi.org/10.1093/cercor/bhad282>
- Jackson, J., Rich, A. N., Williams, M. A., & Woolgar, A. (2017). Feature-selective Attention in Frontoparietal Cortex: Multivoxel Codes Adjust to Prioritize Task-relevant Information. *Journal of Cognitive Neuroscience*, 29(2), 310–321. https://doi.org/10.1162/jocn_a_01039
- Jehee, J. F. M., Brady, D. K., & Tong, F. (2011). Attention Improves Encoding of Task-Relevant Features in the Human Visual Cortex. *Journal of Neuroscience*, 31(22), 8210–8219. <https://doi.org/10.1523/JNEUROSCI.6153-09.2011>
- Jenkinson, M., Bannister, P., Brady, M., & Smith, S. (2002). Improved optimization for the robust and accurate linear registration and motion correction of brain images. *NeuroImage*, 17(2), 825–841. [https://doi.org/10.1016/S1053-8119\(02\)91132-8](https://doi.org/10.1016/S1053-8119(02)91132-8)
- Jenkinson, M., Beckmann, C. F., Behrens, T. E. J., Woolrich, M. W., & Smith, S. M. (2012). Fsl. *NeuroImage*, 62(2), 782–790. <https://doi.org/10.1016/j.neuroimage.2011.09.015>
- Jenkinson, M., & Smith, S. (2001). A global optimisation method for robust affine registration of brain images. *Medical Image Analysis*, 5(2), 143–156. <https://doi.org/10.1016/S1361->

8415(01)00036-6

- Jerde, T. A., & Curtis, C. E. (2013). Maps of space in human frontoparietal cortex. *Journal of Physiology Paris*, 107(6), 510–516. <https://doi.org/10.1016/j.jphysparis.2013.04.002>
- 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. <https://doi.org/10.1016/J.NEURON.2007.02.015>
- Liu, T., Larsson, J., & Carrasco, M. (2007). Feature-based attention modulates orientation-selective responses in human visual cortex. *Neuron*, 55(2), 313–323. <https://doi.org/10.1016/j.neuron.2007.06.030>
- Liu, T., Slotnick, S. D., Serences, J. T., & Yantis, S. (2003). Cortical mechanisms of feature-based attentional control. *Cerebral Cortex (New York, N.Y.: 1991)*, 13(12), 1334–1343. <https://doi.org/10.1093/cercor/bhg080>
- Livingston, K. R., Andrews, J. K., & Harnad, S. (1998). Categorical perception effects induced by category learning. *Journal of Experimental Psychology. Learning, Memory, and Cognition*, 24(3), 732–753. <https://doi.org/10.1037//0278-7393.24.3.732>
- Martinez-Trujillo, J. C., & Treue, S. (2004). Feature-based attention increases the selectivity of population responses in primate visual cortex. *Current Biology: CB*, 14(9), 744–751. <https://doi.org/10.1016/j.cub.2004.04.028>
- McAdams, C. J., & Maunsell, J. H. (1999). Effects of attention on orientation-tuning functions of single neurons in macaque cortical area V4. *The Journal of Neuroscience: The Official Journal of the Society for Neuroscience*, 19(1), 431–441. <https://doi.org/10.1523/JNEUROSCI.19-01-00431.1999>
- McKee, J. L., Riesenhuber, M., Miller, E. K., & Freedman, D. J. (2014). Task Dependence of Visual and Category Representations in Prefrontal and Inferior Temporal Cortices. *Journal of Neuroscience*, 34(48), 16065–16075. <https://doi.org/10.1523/JNEUROSCI.1660-14.2014>

- Meyers, E. M., Freedman, D. J., Kreiman, G., Miller, E. K., & Poggio, T. (2008). Dynamic Population Coding of Category Information in Inferior Temporal and Prefrontal Cortex. *Journal of Neurophysiology*, *100*(3), 1407–1419. <https://doi.org/10.1152/jn.90248.2008>
- Mimica, B., Tombaz, T., Battistin, C., Fuglstad, J. G., Dunn, B. A., & Whitlock, J. R. (2023). Behavioral decomposition reveals rich encoding structure employed across neocortex in rats. *Nature Communications* *2023 14:1*, *14*(1), 1–20. <https://doi.org/10.1038/s41467-023-39520-3>
- Mirabella, G., Bertini, G., Samengo, I., Kilavik, B. E., Frilli, D., Della Libera, C., & Chelazzi, L. (2007). Neurons in Area V4 of the Macaque Translate Attended Visual Features into Behaviorally Relevant Categories. *Neuron*, *54*(2), 303–318. <https://doi.org/10.1016/j.neuron.2007.04.007>
- Navalpakkam, V., & Itti, L. (2007). Search goal tunes visual features optimally. *Neuron*, *53*(4), 605–617. <https://doi.org/10.1016/j.neuron.2007.01.018>
- Newell, F. N., & Bühlhoff, H. H. (2002). Categorical perception of familiar objects. *Cognition*, *85*(2), 113–143. [https://doi.org/10.1016/S0010-0277\(02\)00104-X](https://doi.org/10.1016/S0010-0277(02)00104-X)
- Niell, C. M., & Stryker, M. P. (2010). Modulation of Visual Responses by Behavioral State in Mouse Visual Cortex. *Neuron*, *65*(4), 472–479. <https://doi.org/10.1016/J.NEURON.2010.01.033>
- Oliva, A., & Torralba, A. (2001). Modeling the Shape of the Scene: A Holistic Representation of the Spatial Envelope *. In *International Journal of Computer Vision* (Vol. 42, Issue 3, pp. 145–175).
- Op de Beeck, H. P., Baker, C. I., DiCarlo, J. J., & Kanwisher, N. G. (2006). Discrimination Training Alters Object Representations in Human Extrastriate Cortex. *Journal of Neuroscience*, *26*(50), 13025–13036. <https://doi.org/10.1523/JNEUROSCI.2481-06.2006>
- Op de Beeck, H. P., Wagemans, J., & Vogels, R. (2001). Inferotemporal neurons represent low-dimensional configurations of parameterized shapes. *Nature Neuroscience*, *4*(12),

1244–1252. <https://doi.org/10.1038/nn767>

Pasupathy, A., & Connor, C. E. (1999). Responses to contour features in macaque area V4.

Journal of Neurophysiology, 82(5), 2490–2502.

<https://doi.org/10.1152/jn.1999.82.5.2490>

Rosch, E., Mervis, C. B., Gray, W. D., Johnson, D. M., & Boyes-Braem, P. (1976). Basic

Objects in Natural Categories. In *COGNITIVE PSYCHOLOGY* (Vol. 8, pp. 382–439).

Roy, J. E., Riesenhuber, M., Poggio, T., & Miller, E. K. (2010). Prefrontal Cortex Activity during

Flexible Categorization. *Journal of Neuroscience*, 30(25), 8519–8528.

<https://doi.org/10.1523/JNEUROSCI.4837-09.2010>

Saenz, M., & Boynton, G. M. (2003). The role of competing stimuli in feature-based attention.

Journal of Vision, 3(9), 727. <https://doi.org/10.1167/3.9.727>

Scolari, M., Byers, A., & Serences, J. T. (2012). Optimal Deployment of Attentional Gain during

Fine Discriminations. *Journal of Neuroscience*, 32(22), 7723–7733.

<https://doi.org/10.1523/JNEUROSCI.5558-11.2012>

Scolari, M., & Serences, J. T. (2009). Adaptive allocation of attentional gain. *The Journal of*

Neuroscience, 29(38), 11933–11942. <https://doi.org/10.1523/JNEUROSCI.5642-08.2009>

Serences, J. T., & Boynton, G. M. (2007). Feature-Based Attentional Modulations in the

Absence of Direct Visual Stimulation. *Neuron*, 55(2), 301–312.

<https://doi.org/10.1016/j.neuron.2007.06.015>

Sereno, M., Dale, A., Reppas, J., Kwong, K., Belliveau, J., Brady, T., Rosen, B., & Tootell, R.

(1995). Borders of multiple visual areas in humans revealed by functional magnetic resonance imaging. *Science*, 268(5212), 889–893.

<https://doi.org/10.1126/science.7754376>

Sigala, N., & Logothetis, N. K. (2002). Visual categorization shapes feature selectivity in the

primate temporal cortex. *Nature*, 415(6869), 318–320. <https://doi.org/10.1038/415318a>

Smith, S. M. (2002). Fast robust automated brain extraction. *Human Brain Mapping*, 17(3), 143–

155. <https://doi.org/10.1002/hbm.10062>

Stringer, C., Pachitariu, M., Steinmetz, N., Reddy, C. B., Carandini, M., & Harris, K. D. (2019). Spontaneous behaviors drive multidimensional, brainwide activity. *Science*, *364*(6437). https://doi.org/10.1126/SCIENCE.AAV7893/SUPPL_FILE/AAV7893_STRINGER_SM.PDF

Swisher, J. D., Halko, M. A., Merabet, L. B., McMains, S. A., & Somers, D. C. (2007). Visual Topography of Human Intraparietal Sulcus. *Journal of Neuroscience*, *27*(20), 5326–5337. <https://doi.org/10.1523/JNEUROSCI.0991-07.2007>

Tanaka, K. (1993). Neuronal mechanisms of object recognition. *Science (New York, N.Y.)*, *262*(5134), 685–688.

Tanaka, K. (1996). Inferotemporal cortex and object vision. *Annual Review of Neuroscience*, *19*, 109–139. <https://doi.org/10.1146/annurev.ne.19.030196.000545>

Thorat, S., & Peelen, M. V. (2022). Body shape as a visual feature: Evidence from spatially-global attentional modulation in human visual cortex. *NeuroImage*, *255*, 119207. <https://doi.org/10.1016/j.neuroimage.2022.119207>

Treue, S., & Maunsell, J. H. (1996). Attentional modulation of visual motion processing in cortical areas MT and MST. *Nature*, *382*(6591), 539–541. <https://doi.org/10.1038/382539a0>

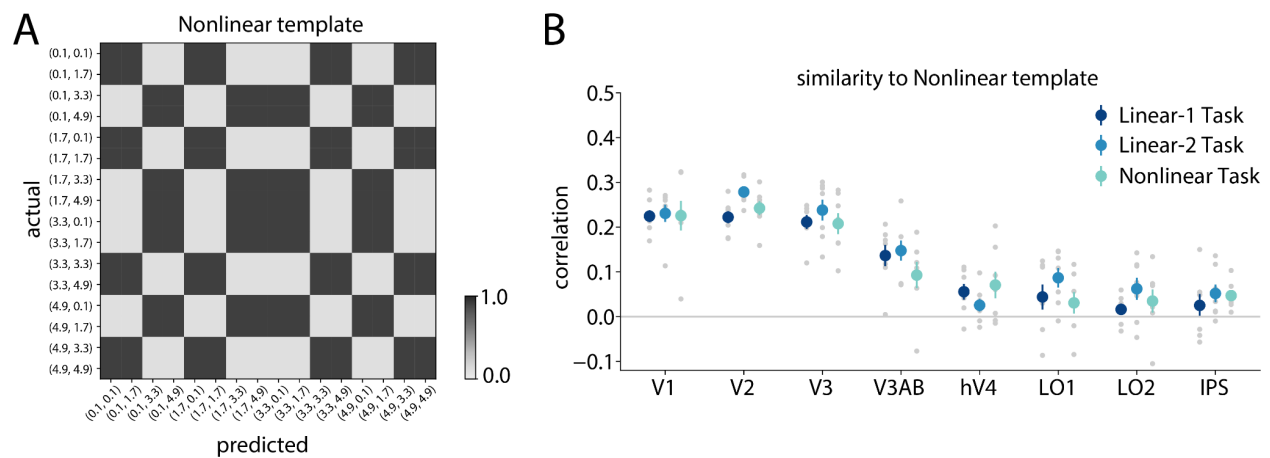
Treue, S., & Maunsell, J. H. (1999). Effects of attention on the processing of motion in macaque middle temporal and medial superior temporal visual cortical areas. *The Journal of Neuroscience: The Official Journal of the Society for Neuroscience*, *19*(17), 7591–7602. <https://doi.org/10.1523/JNEUROSCI.19-17-07591.1999>

Vinberg, J., & Grill-Spector, K. (2008). Representation of Shapes, Edges, and Surfaces Across Multiple Cues in the Human Visual Cortex. *Journal of Neurophysiology*, *99*(3), 1380–1393. <https://doi.org/10.1152/jn.01223.2007>

Wandell, B. A., Dumoulin, S. O., & Brewer, A. A. (2007). Visual field maps in human cortex.

- Neuron*, 56(2), 366–383. <https://doi.org/10.1016/j.neuron.2007.10.012>
- Winawer, J., & Witthoft, N. (2015). Human V4 and ventral occipital retinotopic maps. *Visual Neuroscience*, 32, E020. <https://doi.org/10.1017/S0952523815000176>
- Witkowski, P. P., & Geng, J. J. (2022). Attentional priority is determined by predicted feature distributions. *Journal of Experimental Psychology: Human Perception and Performance*, 48(11), 1201–1212. <https://doi.org/10.1037/xhp0001041>
- Woolrich, M. W., Ripley, B. D., Brady, M., & Smith, S. M. (2001). Temporal autocorrelation in univariate linear modeling of FMRI data. *NeuroImage*, 14(6), 1370–1386. <https://doi.org/10.1006/nimg.2001.0931>
- Yoo, S.-A., Martinez-Trujillo, J. C., Treue, S., Tsotsos, J. K., & Fallah, M. (2022). Attention to visual motion suppresses neuronal and behavioral sensitivity in nearby feature space. *BMC Biology*, 20(1), 220. <https://doi.org/10.1186/s12915-022-01428-7>
- Zahn, C. T., & Roskies, R. Z. (1972). Fourier descriptors for plane closed curves. *IEEE Transactions on Computers*. <https://doi.org/10.1109/TC.1972.5008949>

Supplementary Material



Supplementary Figure 1. Classifier confusion matrix alignment with the *Nonlinear* template does not differ significantly across task conditions. **(A)** Template matrix for the *Nonlinear* task, representing the pattern of similarity expected for a perfect binary representation of the *Nonlinear* categorization scheme. **(B)** The similarity (Pearson correlation coefficient) between the *Nonlinear* template and the actual confusion matrix for each task and ROI. Gray dots represent individual participants, colored circles and error bars represent the mean \pm SEM across 7 participants. A two-way repeated measures ANOVA on these similarity values revealed a main effect of ROI but no main effect of task or ROI x task interaction (ROI: $F_{(7,42)} = 45.16$, $p < 0.001$; Task: $F_{(2,12)} = 1.67$, $p = 0.229$; ROI x Task: $F_{(14,84)} = 1.08$, $p = 0.376$).

Supplementary Table 1. Results of three-way repeated-measures ANOVA tests on the classifier confidence values for far, middle, and near trials, with factors of ROI, task and confidence boundary (i.e., comparing *Linear-1* confidence versus *Linear-2* confidence). Classifier confidence values are shown in Figure 5. All p-values were obtained using a permutation test, see *Methods* for details.

Far trials

	F Value	Num DF	Den DF	p
ROI	43.49	7	42	0.0000
Task	0.65	1	6	0.4532
Boundary	22.09	1	6	0.0025
ROI:Task	0.44	7	42	0.8710
ROI:Boundary	12.07	7	42	0.0000
Task:Boundary	0.65	1	6	0.4559
ROI:Task:Boundary	0.71	7	42	0.6624

Middle trials

	F Value	Num DF	Den DF	p
ROI	24.37	7	42	0.0000
Task	3.10	1	6	0.1243
Boundary	21.21	1	6	0.0036
ROI:Task	0.21	7	42	0.9826
ROI:Boundary	6.51	7	42	0.0000
Task:Boundary	12.16	1	6	0.0120
ROI:Task:Boundary	0.54	7	42	0.8048

Near trials

	F Value	Num DF	Den DF	p
ROI	4.85	7	42	0.0006
Task	6.56	1	6	0.0437
Boundary	0.13	1	6	0.7400
ROI:Task	0.43	7	42	0.8951
ROI:Boundary	1.37	7	42	0.2427
Task:Boundary	6.05	1	6	0.0453
ROI:Task:Boundary	2.06	7	42	0.0647

Supplementary Table 2. Results of two-way repeated-measures ANOVA tests on the *Nonlinear* confidence values for far, middle, and near trials, with factors of ROI and task. Classifier confidence values are shown in Figure 6. All p-values were obtained using a permutation test, see *Methods* for details.

Far trials

	F Value	Num DF	Den DF	p
ROI	42.19	7	42	0.0000
Task	0.53	2	12	0.6034
ROI:Task	0.60	14	84	0.8692

Middle trials

	F Value	Num DF	Den DF	p
ROI	22.25	7	42	0.0000
Task	0.74	2	12	0.5073
ROI:Task	0.78	14	84	0.7015

Near trials

	F Value	Num DF	Den DF	p
ROI	2.25	7	42	0.0480
Task	1.35	2	12	0.3123
ROI:Task	0.83	14	84	0.6414