

# **Food for Thought: Food Preferential Regions in the Human Brain**

Nidhi Jain

CMU-CS-22-108

May 2022

Computer Science Department  
School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213

**Thesis Committee:**

Leila Wehbe, Chair  
Michael J. Tarr

*Submitted in partial fulfillment of the requirements  
for the degree of Masters of Science.*

**Keywords:** neuroscience, statistics, food, encoding models, decoding models, fMRI

## **Abstract**

Existing research has demonstrated functional selectivity in the brain for high level categories of faces, places, bodies, and sensory and motor processes. Food, despite its abundance and importance, has not been considered as a category for which there is a visual selective region, perhaps because of its lack of visual coherence. In this paper, we investigate responsiveness to food in a large scale natural setting via several statistical methods performed on high-resolution fMRI response dataset to natural scenes. We identify two regions consistent across all participants in the high level visual cortex that appear to be functionally selective for food.



## **Acknowledgments**

I would like to thank my advisor Dr. Leila Wehbe for her invaluable guidance the last few years and the fascinating introduction to the research world, Dr. Mike Tarr for his insights, advice, and food puns, Aria Wang for her support and ideas, Maggie Henderson for her constant feedback and willingness to help, and the rest of the invaluable team involved in this project - Andrew Luo, Jacob Prince, and Ruogu Lin. I feel blessed to have had the opportunity to learn from and collaborate with such inspirational and kind scientists. I am also extremely grateful for current and past members of the BrainML lab for the inspiration, motivation, and feedback - Jennifer Williams, Mariya Toneva, Anand Bollu, Catherine Cheng, Zachary Nowak, Tara Pirnia, Joel Ye, and Yuchen Zhou. Finally, I thank my friends and family for their unconditional support.



# Contents

- 1 Introduction** **1**
  
- 2 Methods** **3**
  - 2.1 Dataset . . . . . 3
  - 2.2 Labeling . . . . . 3
  - 2.3 Encoding and Decoding models . . . . . 5
  - 2.4 Investigating Intra-Region Patterns . . . . . 5
  - 2.5 Embedding-based Image Clustering . . . . . 5
  
- 3 Results** **7**
  - 3.1 Encoding Models . . . . . 7
  - 3.2 Decoding Models . . . . . 7
  - 3.3 Apparent Intra-Region Semantic Axes . . . . . 9
  - 3.4 Generalizing Patterns to All Food-related Stimuli . . . . . 9
  - 3.5 Voxel Embedding Clusters vs Image Feature Embedding Clusters . . . . . 12
  
- 4 Discussion** **15**
  
- Bibliography** **17**





# List of Figures

- 2.1 The 1000 images viewed by all 8 subjects in NSD were manually relabeled in order to investigate responsiveness to natural food images. The top half of this figure shows example images that were labeled as well as their corresponding labels, while the bottom half demonstrates the organizational structure used when labeling a given image. Each of the images were given at least one label within each of the three categories of location, content, and image perspective. Multiple content labels could be and often were used for a given image . . . . . 4
  
- 3.1 This figure shows voxels that are determined to be food-related based off of two differing methods. The top image shows significant voxels (white) from a 1-sided t-test comparing food weights from a trained OLS model against all other non-food object weights for all 8 subjects. The middle image shows classification accuracy for voxel-wise searchlight decoding for 4 subjects, with darker blue voxels signifying higher accuracy. The bottom image shows significant voxels from a 1-sided t-test comparing food weights against face weights (red) as well as the results for a 1-sided t test comparing face weights against food weights (blue) for 4 subjects. It is notable here to compare significant 'food' voxels to 'face' voxels because of the high number of face images in our dataset, as well as the robust reliability of the FFA. The food-related regions resulting from both methods overall align with each other, suggesting that these regions are reliably more responsive to food stimuli than other categories. . . . . 8

3.2	Applying PCA on the proposed food region’s responses provides more insight into the structure of these responses. Above, we can see the different regions corresponding to different components, as well as the top images for each component averaged across all subjects. The top images display each components projection onto the brain, demonstrating the generally positively (green) and negatively (brown) contributing voxels. On the top right, we see in red an example of the region that we have identified as the relevant food region for a given subject, generated based off of the intersection of the food vs face encoding model test and manually identified visual cortex regions. In the graphs, we see stimuli that most contribute to various components, as well as their placement. We see that each component overall highlights different semantic characteristics of the visual stimuli. The first component appears captures the surface area of food in the image, the second distinguishes places, and the third seems to capture a combination of spatial frequency and people. These semantic categories are able to explain a large component of the variance resulting from these food regions. . . .	10
3.3	Analysis on predicting voxel activity from unseen images points to the same food-selective regions as previously proposed. We observe 4 COCO label weights corresponding to each voxel following a voxel-wise ridge regression on COCO labels (left). We also observe the voxel-wise improved $R^2$ values from including food labels when predicting voxel activity (right). Both the weights corresponding to individual food labels as well as the overall impact of the aggregate food labels highlight our proposed food regions. These results also point to the conclusion that this region appears food selective. . . . .	11
3.4	We observe the highest voxel-activity inducing images in each cluster. This helps us identify how our proposed food region reacts to different images. On the left we see in red the voxels identified as the proposed food region, and thus used here to cluster the images. In the middle we see top images corresponding to our first cluster, and on the right we see the top images corresponding to the second cluster. We observe that this region appears to respond differently to food than to other images. . . . .	12
3.5	We observe image clusters based on voxel-response embeddings (top). We then compare this to image clusters based on CLIP embeddings (bottom). Two clusters emerge from voxel-response embeddings, while 4 main clusters emerge from CLIP. The two clustering patterns show little to no correlation with each other, demonstrating that this proposed food region involves more complex, top down processing rather than solely low level visual and semantic features. . . . .	14

# Chapter 1

## Introduction

The representation of high-level visual information in the human brain has been marked by the phenomenon of selectivity for visual categories of high ecological importance. There are multiple brain regions that show preferential neural responses to such categories: faces, bodies, and places [1, 2, 3, 4]. Independent of any particular theory on the origins and specificity of these functional brain regions [5, 6], we can agree that they exist because they instantiate processes and representations for categories that are highly relevant for day-to-day behavior. In a similar vein, food is a category that is evolutionarily relevant – finding nourishment being more ancient than social interaction and, arguably, more fundamental to survival. It is surprising that food has not been identified as a visual category for which localized preferential neural responses are observed. Why is visual selectivity for food not low hanging fruit?

While it is known that the visual presentation of food images prompts a range of brain responses, including both cognitive and emotional effects, the concurrence between studies using food images has been quite low [7]. Only 41% of the studies in a meta analysis contributed to food-related clusters in the fusiform gyrus (bilaterally) and orbitalfrontal cortex (left) [7]. And in cases where food-related responses have been observed, they typically have been attributed to increased attention to food images because of participants’ mental states and/or physiological factors [7, 8]. As such, there have been no clear tests as to whether food elicits preferential responses in the human visual system.

An important factor that may influence prior studies of food-elicited neural responses is context. One of the characteristics of category-selective neural responses is that they are elicited automatically by the category in question [9]. In contrast, as already mentioned, in many of the extant studies using food images, food-selective responses appear to be elicited only under certain affective states or physiological conditions [7]. We posit that this apparent inconsistency in detecting food-preferential neural responses is, in part, a result of relying on isolated food and non-food “posed” visual stimuli, and that many studies have failed to consider food in environmentally relevant contexts. Decontextualized food images – which vary widely in visual appearance – may not be sufficiently salient or realistic enough so as to elicit food-related processing by default (although such images may do so with task-driven attention). In contrast, both faces and bodies exhibit very little visual variation within category, while places are typ-

ically tested using natural scene images. In a nutshell, the extreme within-category variability in food appearance may contribute in several ways that may render identifying food-preferential brain regions more challenging than other ecologically important categories. First, without context, *images* of food may not present as food *qua* food to the visual system. That is, category membership for images of food may be ambiguous absent the presence of associated semantic information that helps, possibly in a top-down manner, constrain and narrow the categorical content of the image. Second, because of the high visual variability for food as a category, detecting significant food-driven responses may require more sensitive designs than provided by standard neuroimaging designs (which typically rely on small numbers of images and trials per condition and therefore depend on high within-condition similarity).

Our study addresses these issues in two ways. First, real-world images, drawn from the the COCO dataset [10], were used for both the food and non-food conditions. Second, functional MRI (fMRI) data in response to viewing these images was collected on a massive scale [11], thereby improving our ability to detect effects across conditions. To preview our most important result, we reliably identify two distinct regions in high-level visual cortex that are preferentially responsive to food images. As with other robust category preferences, we take our results to indicate that the human visual system instantiates (through unknown mechanisms) processes and representations that support perceiving and reasoning about food. In that food is incontrovertibly an ecologically critical category, we view this finding as highly consistent with earlier findings of preferential mechanisms for the perception of faces, bodies, and places.

# Chapter 2

## Methods

### 2.1 Dataset

To investigate responsiveness to food in a large scale natural setting, we take advantage of the Natural Scenes Dataset (NSD), a dataset of high-resolution fMRI responses to natural scenes via a recognition task [11]. These natural scene images are pulled from the annotated Microsoft Common Objects in Context (COCO) dataset [10]. NSD has brain response data from 8 screened subjects that each see on average 9500 natural scene images over the course of a year. These viewings were administered during 30-40 scan sessions. Of the 70,566 total unique images viewed across participants, 1000 are viewed by all 8 participants. We took special interest in these 1000 images' response data when conducting analysis in order to ensure consistency in results.

### 2.2 Labeling

To extract more unlabeled information and ensure full reliability of the 1000 shared images that are viewed by all subjects, we methodologically relabel them based on 3 main categories. We use a categorical hierarchical structure as shown in Figure 2.1 to label these images with their location classification, object classification (including the binary existence of food), and perspective classification of the image. Image perspective is discretized into *zoom*, *reach* or *large-scale*. *Zoom* signifies an apparent focusing by the camera lens, thereby likely concentrated on one object and excluding surrounding information. *Reach* images demonstrate affordances by displaying objects at a human-reachable distance [12]. *Large-scale* images encompass the remaining images, which include an image of a general scene as opposed to one or more close up objects. The image perspective category's vague nature leaves it vulnerable to variation in labeling. To avoid this variation and ensure consistency, we undergo several rounds of labeling and verification.

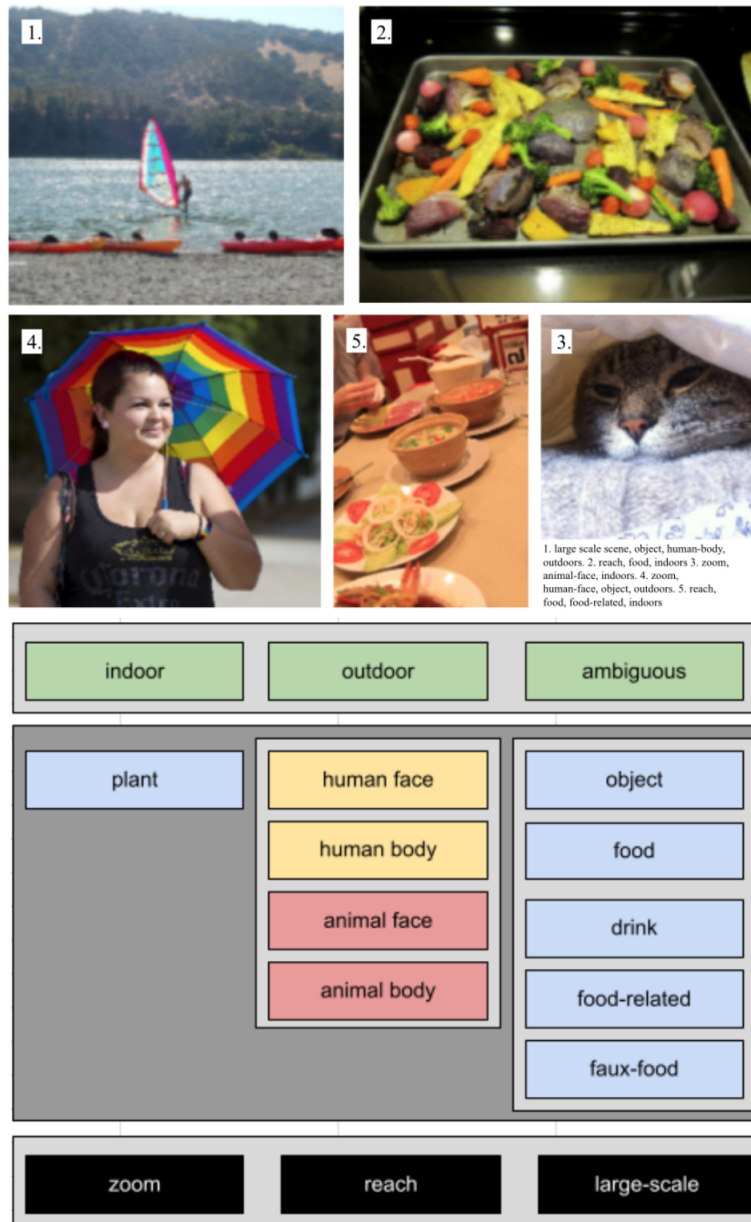


Figure 2.1: The 1000 images viewed by all 8 subjects in NSD were manually relabeled in order to investigate responsiveness to natural food images. The top half of this figure shows example images that were labeled as well as their corresponding labels, while the bottom half demonstrates the organizational structure used when labeling a given image. Each of the images were given at least one label within each of the three categories of location, content, and image perspective. Multiple content labels could be and often were used for a given image

## 2.3 Encoding and Decoding models

Using these hand-labeled images, we can reliably perform statistical methods to identify voxels especially responsive to given categories. Encoding all 16 labels into a single binary vector per image, we utilize voxel-wise ordinary least squares (OLS) encoding models to predict each individual voxel response to a given stimuli. Identifying voxels more responsive to category  $A$  over category  $B$  then involves a 1 sided  $t$ -test between the respective learnt model coefficients for each of the two categories. We use these methods to identify voxels that are more responsive to food than other categories.

While an encoding model is able to provide some insight into single-voxel selectivity through response predictions, a decoding model can uncover distributed pattern-level representations of visual features. Taking advantage of this characteristic of decoding models, we run voxel-wise searchlight classifiers for the existence of food and further inspect high-performing regions. Regions emerging from this method thus present processing information via a pattern that our model is able to exploit.

## 2.4 Investigating Intra-Region Patterns

Encoding and decoding models help identify proposed regions especially responsive to food, but are unable to provide further insight regarding the structure of these 'food regions'. We run a principal component analysis to better understand possible structure and/or correspondence in these food regions.

At this point, our methods have mostly focused on identifying apparent regions through analysis of responses to the 1000 shared images. Using rich, specific COCO annotations (including specific types of food), we can further investigate the remaining non-hand labeled images to verify proposed food-selective regions. We use a ridge regression model on COCO labels and observe the resulting voxel-wise weights for a specific label. This method allows us to identify regions that appear to have higher correlations to a given label.

## 2.5 Embedding-based Image Clustering

In addition to investigating the food category at a fine grained level via label-related activation, we perform image-level response analysis via image clustering. Using the proposed food-selective region as a voxel mask, we extract each image's corresponding food-voxel activity. We then perform K-means clustering on this activity and investigate the resulting groups for possible image-level patterns.

To better investigate visual and semantic patterns in these voxel-based clusters and discuss potential response impact from top-down cortical influences, we also cluster visual and semantic

embeddings without using brain activity of these images and compare the image clustering results to voxel clustering results. To obtain these visual and semantic embeddings we use two deep learning models, CLIP and ResNet-50 [13, 14]. CLIP, trained on both images and text, allows us to extract features with semantic and visual meaning. ResNet-50, training on solely images, results in purely visual feature-based embeddings. Similar clustering results would lead us to believe that responses may significantly weigh similar features as these models, while varying clustering results provide further evidence that deeper processing may be involved than purely low level or semantic features.



# Chapter 3

## Results

### 3.1 Encoding Models

We find two main areas in the visual cortex that appear consistently selective across all subjects for food when compared to other non-food categories via significance tests on OLS encoding model weights, as shown in the top of Figure 3.1. This finding persists across all 8 subjects after performing 1 sided significance tests when comparing food weights to weighted non-food object weights. These two regions appear to respond to food more significantly than other non-food categories. When considering a one sided significance test between 2 labels, for example, food and face, the significant voxels that achieve a  $p$ -value lower than 0.05 are identified as more responsive to food than faces. As in the top portion of Figure 3.1, the food-significant voxels are often extremely close to 0, and thus highly significant, when compared against other categories. We also note that this spatial 'food' pattern persists even when compared against the *reach* label. These voxels highlighted in Figure 3.1 thus seem to be most activated by the existence of food in the stimuli as opposed to other categories, pointing towards an apparent functional selectivity.

As a sanity check of this method, we verify that performing a similar significance test for faces reveals the fusiform face area (FFA), as shown in blue in the bottom portion of Figure 3.1 [2]. We choose the FFA here for our sanity check due to the high occurrence of faces in the dataset and reliability of this region. We notice that the identified food region (red) and face region (blue) via this method appear to have completely independent selections, further highlighting an apparent selectivity for food. Nonetheless, it is clear that these proposed 'food-significant' voxels appear consistently more responsive to food stimuli than other stimuli.

### 3.2 Decoding Models

The high-accuracy voxels from the searchlight decoding model overall align spatially with the significant voxels from the encoding model, as is clear in Figure 3.1. We note that high-accuracy voxels do not necessarily positively correspond to food-selective regions, but rather to any region which helps identify whether or not the given stimuli contains food. So, high-accuracy voxels can signify both a food-selective region and an anti-food selective region that especially helps

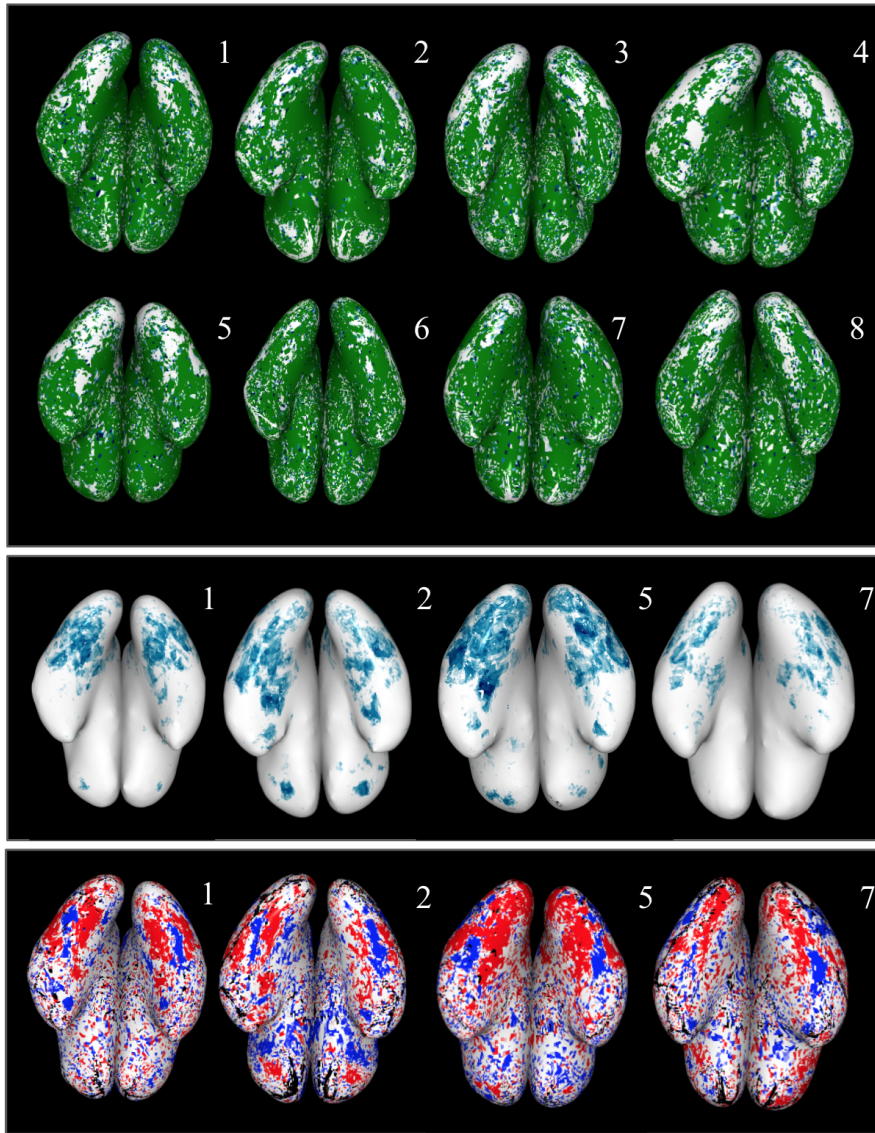


Figure 3.1: This figure shows voxels that are determined to be food-related based off of two differing methods. The top image shows significant voxels (white) from a 1-sided t-test comparing food weights from a trained OLS model against all other non-food object weights for all 8 subjects. The middle image shows classification accuracy for voxel-wise searchlight decoding for 4 subjects, with darker blue voxels signifying higher accuracy. The bottom image shows significant voxels from a 1-sided t-test comparing food weights against face weights (red) as well as the results for a 1-sided t test comparing face weights against food weights (blue) for 4 subjects. It is notable here to compare significant 'food' voxels to 'face' voxels because of the high number of face images in our dataset, as well as the robust reliability of the FFA. The food-related regions resulting from both methods overall align with each other, suggesting that these regions are reliably more responsive to food stimuli than other categories.

determine that the stimuli is not food-related. The high performing regions from this searchlight model are generally consistent across all subjects. Similar regions as those found in the encoding model are also observed to be highly correlated with the stimuli food classification accuracy. The decoding model’s inclusion of the encoding model’s resulting food-selective region once again emphasizes our finding of this food-selective region.

Both the encoding and decoding models confirm that these regions seem to be consistently more responsive to food than other categories. However, these models provide little insight into the more fine grained semantic and spatial structure of the responses within this food region. To discover these intra-region patterns, we run a principled component analysis (PCA) on the matrix of all subjects’ voxels in the selected visual cortex by all shared food images. We then observe each of the principal components’ corresponding voxel values as well as the top contributing images to each component.

### **3.3 Apparent Intra-Region Semantic Axes**

Running PCA on the hypothesized food region leads to a further spatial breakdown of various ‘food’ areas, as is shown in Figure 3.2. These breakdowns via principal components are able to account for a significant portion of the variance of this region’s responses. The top components explain 34.31, 12.68, and 11.16 percent of the variability, respectively. Each of the uncovered components’ axes also seems to correspond to a general semantic pattern of image content. The first component appears to represent the surface area of food, distinguishing images with few pixels corresponding to food from those with mostly food pixels. The second principal component distinguishes images with a focus on food from zoomed out images of places related to food. The third principal component distinguishes images with a focus on food from zoomed out images of people eating food, with social settings being at the end of the spectrum. These results highlight the importance of food as a category that is easily disentangled from the other categories that are important in the fusiform cortex, namely faces and places.

### **3.4 Generalizing Patterns to All Food-related Stimuli**

To further consider potential intra-region semantic breakdowns beyond our PCA results, we also directly inspect voxel response correlations with fine-grained food categories. To increase the size of our dataset and therefore reliability, we observe how voxels react to specific food categories from images not yet inspected when proposing the food region. In order to use images not yet inspected, we consider the remaining images that were not viewed by all 8 subjects. We run a voxel-wise ridge regression on the provided fine-grained COCO labels for these images to predict voxel activity. Then, we visualize the resulting weights per label on each voxel, as shown in Figure 3.3. In order to observe the which voxels’ prediction accuracies benefit most from the addition of food labels, we perform the ridge regression on all labels, as well as all

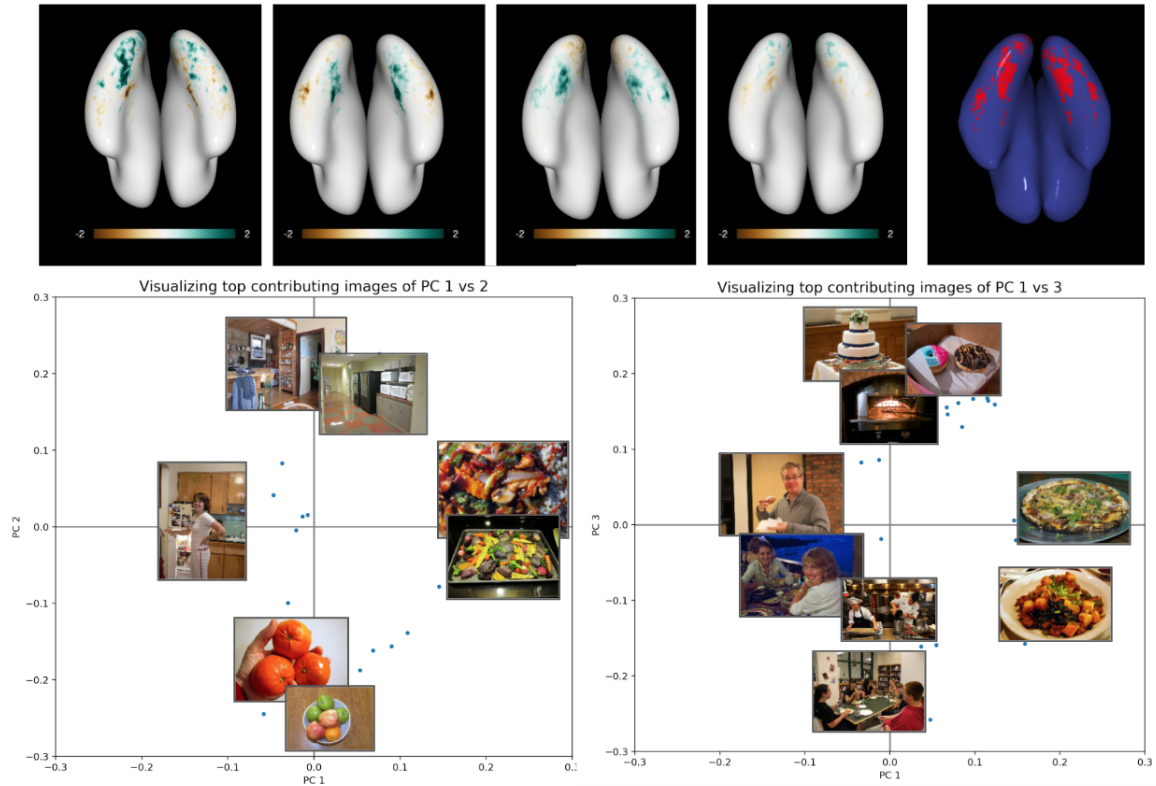


Figure 3.2: Applying PCA on the proposed food region’s responses provides more insight into the structure of these responses. Above, we can see the different regions corresponding to different components, as well as the top images for each component averaged across all subjects. The top images display each components projection onto the brain, demonstrating the generally positively (green) and negatively (brown) contributing voxels. On the top right, we see in red an example of the region that we have identified as the relevant food region for a given subject, generated based off of the intersection of the food vs face encoding model test and manually identified visual cortex regions. In the graphs, we see stimuli that most contribute to various components, as well as their placement. We see that each component overall highlights different semantic characteristics of the visual stimuli. The first component appears captures the surface area of food in the image, the second distinguishes places, and the third seems to capture a combination of spatial frequency and people. These semantic categories are able to explain a large component of the variance resulting from these food regions.

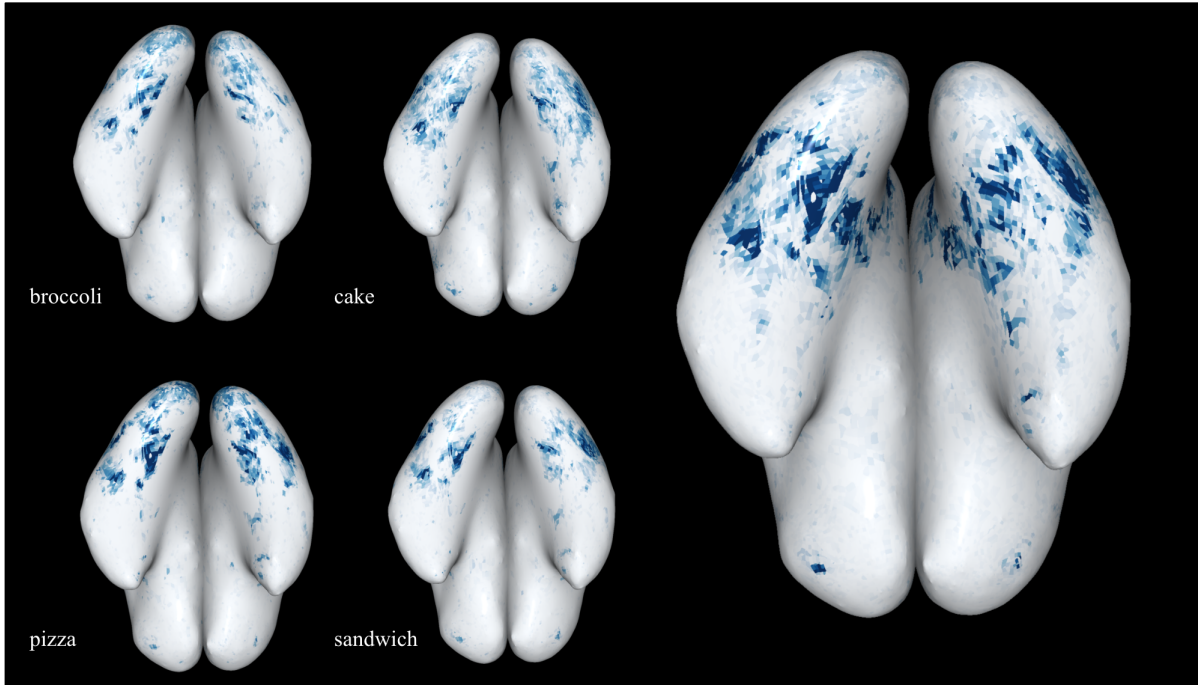


Figure 3.3: Analysis on predicting voxel activity from unseen images points to the same food-selective regions as previously proposed. We observe 4 COCO label weights corresponding to each voxel following a voxel-wise ridge regression on COCO labels (left). We also observe the voxel-wise improved  $R^2$  values from including food labels when predicting voxel activity (right). Both the weights corresponding to individual food labels as well as the overall impact of the aggregate food labels highlight our proposed food regions. These results also point to the conclusion that this region appears food selective.

but food labels. We then compare the  $R^2$  values voxel-wise between results from including all labels during training, and results from including all but food labels in training. We observe one general pattern of voxels whose prediction performance improved due to the inclusion of food-related labels, and the same underlying pattern for high valued voxel-wise weights for individual fine grained food labels (*i.e.* *cake*, *sandwich*). These activations across the fine grained food labels appear consistent within the food category as well as across subjects. This pattern also corresponds strongly with our previously proposed food region. Both the  $R^2$  improvement and the fine-grained labels' corresponding weights suggest that voxels in these regions consistently appear more responsive to food than other cortical areas. The responses re-emphasize our identification of a food-selective region.



Figure 3.4: We observe the highest voxel-activity inducing images in each cluster. This helps us identify how our proposed food region reacts to different images. On the left we see in red the voxels identified as the proposed food region, and thus used here to cluster the images. In the middle we see top images corresponding to our first cluster, and on the right we see the top images corresponding to the second cluster. We observe that this region appears to respond differently to food than to other images.

### 3.5 Voxel Embedding Clusters vs Image Feature Embedding Clusters

In addition to observing potential intra-region patterns in the proposed food region, we also use a larger subset of COCO images to perform image clustering via their respective voxel embeddings in efforts to understand emerging image-level patterns. In this test, our new subset of images also includes an even split of non-food and food images from images that were not used to determine the originally proposed food region. Using the previously proposed mask of voxels identified as food-selective, we create a subset of corresponding voxel responses for each image. We then use K-means to cluster the voxel embeddings corresponding to these images to reveal patterns of responsiveness in this region. 2 clusters emerge, each containing a group of voxel response vectors corresponding to images. We now consider the images corresponding to this voxel activity, and thus the 2 clusters of images. We notice that 1011 of the 1638 food images gather in one cluster, and observe the images with the highest voxel activity in each cluster (Figure 3.4). The top images in each cluster have significant semantic differences, with one cluster clearly corresponding highly with food. These top cluster images further confirm our proposal of this region as selective to food.

We use voxel responses to further cluster food images and observe the resulting clusters. Especially when considering the visual cortex, it is important to isolate causes for variations in brain responses from solely image feature variation. To help understand the extent to which variations in food responses are due to visual characteristics, we use CLIP and ResNet to cluster food-related images and observe resulting patterns. We then observe whether these patterns are consistent with the clusters emerging from voxel-based image clustering. CLIP considers both visual and semantic features when training, while ResNet-50 only considers visual features. In

both outputs, we notice semantic clusters that do not have clear correspondence with our resulting food clusters (Figure 3.5), suggesting likely top down processing that goes further than raw visual features. We also note the possibility of these voxels responding to the composition of visual features. Overall, this lack of consistency in our two clustering results demonstrates that the voxel embeddings do not appear to represent information in the same way as our raw image-feature embeddings.

Finally, it is unlikely that our finding of food selectivity can be solely attributed to greater attention or higher intrinsic visual saliency for food relative to non-food. First, both human faces and bodies would seem to be subject to the same kinds of effects. However, attentional/saliency differences are not the preferred explanation for the basic finding of face or body selectivity [15]. Moreover, within our study, faces and bodies comprised a reasonable proportion of the non-food contrast images, yet food selectivity was robust across these comparison categories. Second, the images used in the fMRI data collection, drawn from the COCO image dataset, all depict complex natural scenes containing consistently high saliency objects and actions. It is therefore doubtful that the non-food images drew less attention or appeared less salient than the food images. If anything, real-world images lead to better control in across-category comparisons because of diminished differences in attention and saliency.

Image clusters emerging from voxel embeddings:



Image clusters emerging from CLIP embeddings:

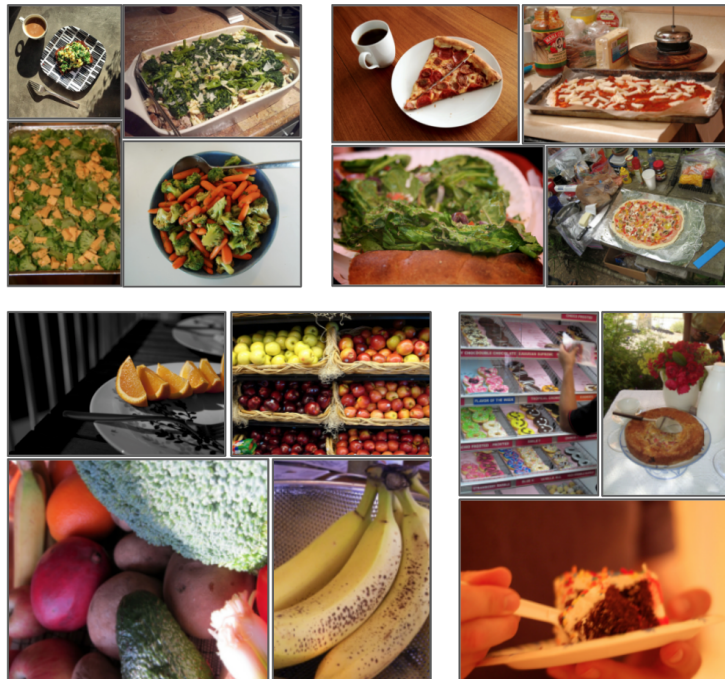


Figure 3.5: We observe image clusters based on voxel-response embeddings (top). We then compare this to image clusters based on CLIP embeddings (bottom). Two clusters emerge from voxel-response embeddings, while 4 main clusters emerge from CLIP. The two clustering patterns show little to no correlation with each other, demonstrating that this proposed food region involves more complex, top down processing rather than solely low level visual and semantic features.



# Chapter 4

## Discussion

How are knowledge representations organized in the human brain? Within the visual system, one of the hallmarks of the past several decades has been *category selectivity* for faces, bodies, and places [1, 2, 3, 4]. Consistent with the ecological importance of these categories, we predicted and found selectivity within the visual system for food. However, in contrast to other domains that show selectivity, food responses are much less localized, suggesting a complex network of brain regions underlying food preference. One explanation for this pattern of results may be that, relative to faces, bodies, and places, food appearances are quite variable. As such, it is unlikely that there is a set of lower-level visual features or high-level shape structures that consistently correspond to food (in contrast, see [16, 17]).

How then, do food preferential mechanisms and representations arise in the human brain? Similar to human language, domain-relevant perceptual inputs can vary widely depending on the cultural and physical environment. Consequently, it is high-level semantic and structural biases (and not low-level features) that are critical for learning the domain. Under this view, learned representations are only loosely constrained at the surface level, but still reflect common underlying mechanisms that have emerged over the course of evolution due to selection for general learning abilities (the “Baldwin Effect” [18, 19]). That is, because the mechanisms supporting domain acquisition are able to flexibly respond to variations in inputs, these learning capacities are preserved across evolution. In this light, as a consistent property of knowledge organization, food selectivity is likely to have emerged as neural preference shaped heavily by top-down knowledge and semantic associations rather than low-level properties of the inputs themselves.



# Bibliography

- [1] J Sergent, S Ohta, and B MacDonald. Functional neuroanatomy of face and object processing: A positron emission tomography study. *Brain*, 115:15–36, 1992. 1, 4
- [2] N Kanwisher, J McDermott, and M M Chun. The fusiform face area: a module in human extrastriate cortex specialized for face perception. *J Neurosci*, 17(11):4302–4311, 1997. 1, 3.1, 4
- [3] R Epstein and N Kanwisher. A cortical representation of the local visual environment. *Nature*, 392(6676):598–601, 1998. doi: 10.1038/33402. 1, 4
- [4] P E Downing, Y Jiang, M Shuman, and N Kanwisher. A cortical area selective for visual processing of the human body. *Science*, 293(5539):2470–2473, 2001. doi: 10.1126/science.1063414. 1, 4
- [5] N Kanwisher. Domain specificity in face perception. *Nat. Neurosci.*, 3(8):759–763, 2000. 1
- [6] M J Tarr and I Gauthier. FFA: a flexible fusiform area for subordinate-level visual processing automatized by expertise. *Nat Neurosci*, 3(8):764–769, 2000. doi: 10.1038/77666. 1
- [7] VanderLaan. The first taste is always with the eyes: A meta-analysis on the neural correlates of processing visual food cues. *NeuroImage*, 55(1):296–303, 2011. doi: 10.1016/j.neuroimage.2010.11.055. 1
- [8] Ruud van den Bos and Denise de Ridder. Evolved to satisfy our immediate needs: Self-control and the rewarding properties of food. *Appetite*, 47(1):24–29, 2006. 1
- [9] I Gauthier, M J Tarr, A W Anderson, P Skudlarski, and J C Gore. Activation of the middle fusiform 'face area' increases with expertise in recognizing novel objects. *Nat Neurosci*, 2(6):568–573, 1999. doi: 10.1038/9224. 1
- [10] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision – ECCV 2014*, pages 740–755. Springer International Publishing, 2014. 1, 2.1
- [11] Allen. A massive 7t fmri dataset to bridge cognitive neuroscience and artificial intelligence. *Nat Neurosci*, 25(1):116–126, 2022. doi: 10.1038/s41593-021-00962. 1, 2.1
- [12] Emilie L Josephs, Haoyun Zhao, and Talia Konkle. The world within reach: an image database of reach-relevant environments. *Journal of Vision*, 21(7):14–14, 2021. 2.2

- [13] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 2.5
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 2.5
- [15] Maura L. Furey, Topi Tanskanen, Michael S. Beauchamp, Sari Avikainen, Kimmo Uutela, Riitta Hari, and James V. Haxby. Dissociation of face-selective cortical responses by attention. *Proceedings of the National Academy of Sciences*, 103(4):1065–1070, 2006. doi: 10.1073/pnas.0510124103. 3.5
- [16] Shahin Nasr and Roger B H Tootell. A cardinal orientation bias in scene-selective visual cortex. *J Neurosci*, 32(43):14921–14926, 2012. doi: 10.1523/JNEUROSCI.2036-12.2012. 4
- [17] Xiaomin Yue, Irene S Pourladian, Roger B H Tootell, and Leslie G Ungerleider. Curvature-processing network in macaque visual cortex. *Proc Natl Acad Sci U S A*, 111(33):E3467–75, 2014. doi: 10.1073/pnas.1412616111. 4
- [18] J. Mark Baldwin. A new factor in evolution. *The American Naturalist*, 30(354):441–451, 1896. ISSN 00030147, 15375323. URL <http://www.jstor.org/stable/2453130>. 4
- [19] Patrick Bateson. The Active Role of Behaviour in Evolution. *Biology and Philosophy*, 19(2):283–298, 2004. ISSN 1572-8404. doi: 10.1023/B:BIPH.0000024468.12161.83. URL <https://doi.org/10.1023/B:BIPH.0000024468.12161.83>. 4
- [20] Nancy Kanwisher. Functional specificity in the human brain: A window into the functional architecture of the mind. *Proceedings of the National Academy of Sciences*, 107(25):11163–11170, 2010.
- [21] Michel Goossens, Frank Mittelbach, and Alexander Samarin. *The L<sup>A</sup>T<sub>E</sub>X Companion*. Addison-Wesley, Reading, Massachusetts, 1993.
- [22] Donald Knuth. Knuth: Computers and typesetting. URL <http://www-cs-faculty.stanford.edu/~dunk/abcde.html>.
- [23] Raymond Noble and Denis Noble. Was the watchmaker blind? or was she one-eyed? *Biology*, 6(4), 2017. ISSN 2079-7737. doi: 10.3390/biology6040047. URL <https://www.mdpi.com/2079-7737/6/4/47>.
- [24] William D. S. Killgore and Deborah A. Yurgelun-Todd. Positive affect modulates activity in the visual cortex to images of high calorie foods. *International Journal of Neuroscience*, 117(5):643–653, 2007. doi: 10.1080/00207450600773848.
- [25] Carolijn Ouwehand and Esther K Papiés. Eat it or beat it. the differential effects of food temptations on overweight and normal-weight restrained eaters. *Appetite*, 55(1):56–60, 2010.