
AlignedCut: Visual Concepts Discovery on Brain-Guided Universal Feature Space

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study the intriguing connection between visual data, deep networks, and the
2 brain. Our method creates a universal channel alignment by using brain voxel
3 fMRI response prediction as the training objective. We discover that deep net-
4 works, trained with different objectives, share common feature channels across
5 various models. These channels can be clustered into recurring sets, correspond-
6 ing to distinct brain regions, indicating the formation of visual concepts. Tracing
7 the clusters of channel responses onto the images, we see semantically meaning-
8 ful object segments emerge, even without any supervised decoder. Furthermore,
9 the universal feature alignment and the clustering of channels produce a picture
10 and quantification of how visual information is processed through the different
11 network layers, which produces precise comparisons between the networks.

12 1 Introduction

13 Introducing a novel approach, Yang et al. (2024) has successfully established a method of computing
14 a mapping between the brain and deep-nets, effectively linking two black boxes. The brain fMRI
15 prediction task allows for visualizing information flow from layer to layer, using the brain as an
16 analysis tool.

17 If a picture is worth a thousand words, the
18 main idea is that the brain’s thousands of voxels
19 can be thought of as alphabets for these words
20 that describe an image. Just as alphabets must
21 be combined to form words and phrases with
22 meanings, we need to find the grouping of brain
23 voxels and their network channel counterparts
24 to understand their meaning (Figure 1).

25 Our main discovery is that while the network
26 layer structure differs, channel feature corre-
27 spondence exists across networks with a shared
28 encoding of reoccurring visual concepts. This paper builds upon the idea of ‘Rosetta stone’ neurons
29 (Dravid et al., 2023), which find channels across networks that share similar image responses in bi-
30 nary segmentation. If channels are alphabets, ‘Rosetta stone’ provides an alphabet-level translation
31 between networks.

32 Individual channel-level analysis could miss feature correspondence across networks at finer and
33 coarser levels. On a finer level, because the channels are invariant up to a linear transformation, we
34 might miss a reconstituted feature constructed from a composition of existing channels. On a coarse
35 level, the channels can be combined and clustered to form a bigger ‘Rosetta’ concept.

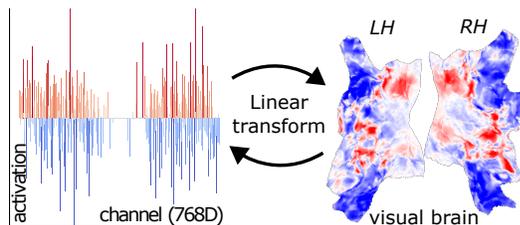


Figure 1: Transform the hidden channel activation of deep-nets into visual brain voxels’ response.

36 To address fine-level channel analysis, we use brain voxel response as a reference signal and linearly
 37 transform channels for each network into a shared space sufficient for brain fMRI prediction. This
 38 process produces a universal feature space that aligns channel features across the layers and models.
 39 To find bigger visual concepts, one can start with Neuroscience knowledge of brain regions (ROIs)
 40 with specific brain functionality, i.e., V1, V4, and EBA. While tracing the mapping of the ROIs to
 41 channels can produce visual concepts (Figure 2), brain regions don't function in isolation.

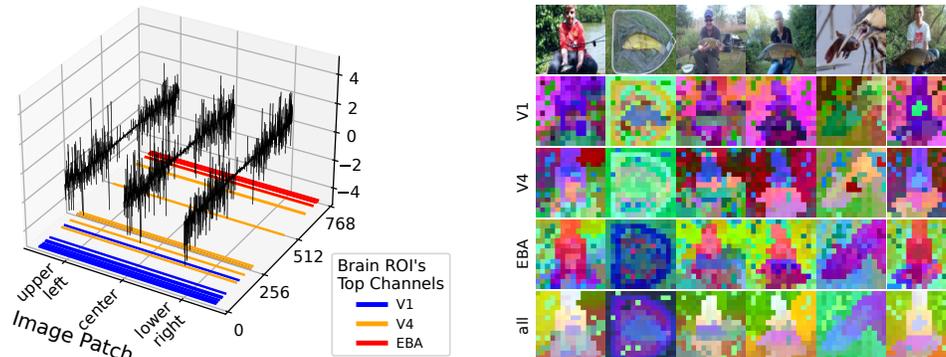


Figure 2: From the 768D feature on CLIP layer-6, we extract different levels of segmentation by restricting the use of a subset of channels. **Left:** Channel activation on example image patches. The ordering of channels is sorted from the early brain to the late brain by their weights for brain voxels. **Right:** Spectral clustering on each subset of channels filtered by each brain ROI (V1, V4, EBA), image pixels colored by 3D spectral-tSNE of top 10 eigenvectors.

42 Instead of searching through all possible channel grouping combinations, our first insight is that
 43 we can create a channel grouping hypothesis by examining channels from each pixel's perspective.
 44 Think of the pixels and channels forming a bipartite graph; each channel produces a per-pixel re-
 45 sponse (image activation map), defining the graph edge between the pixels and channels. Taking the
 46 perspective of pixels, one can collect graph edges incident on each pixel into a vector, which can be
 47 thresholded to produce a hypothesis grouping over channels.

48 Our second insight is that if a channel grouping hypothesis repeats across images, layers, and mod-
 49 els, it is highly unlikely to be accidental and, therefore, signals meaningful visual concepts.

50 We formulate this clustering problem as a graph partition task. The graph nodes are the product
 51 space of pixels and layers. We apply spectral clustering to produce k-top eigenvectors. We take
 52 advantage of two properties of spectral clustering: it makes 1) soft-cluster embedding space in the
 53 form of eigenvectors and 2) hierarchical clustering by varying the number of eigenvectors.

54 We made the following discoveries. First, shared channel sets, reoccurring across layers and models,
 55 predict response in distinct brain regions. By tracing the channel activation to the known brain ROI
 56 properties, we observe that the channel cluster encodes visual concepts at various levels of visual
 57 abstraction.

58 Second, meaningful object segments can emerge by tracing the channel cluster responses onto each
 59 image. We observed that some channel clusters produce figure/ground separation while others pro-
 60 duce fine-grained category classification. Our image segmentation requires no additional segmenta-
 61 tion decoder and uses only a simple distance measure over the eigenvectors.

62 Finally, the universal feature alignment and the spectral clustering of channels produce a picture and
 63 quantification of how visual information is processed through the different network layers.

64 While these discoveries are promising, there are two main technical hurdles to overcome to verify
 65 them on a large scale. Our method rests upon a crucial assumption: the channels across the different
 66 layers and models can be mapped into a shared space. While brain prediction over thousands of
 67 voxels can provide strong guidance for this alignment, an additional constraint would be needed
 68 when the shared space has a large dimension (suitable for expressiveness). We use clustering as
 69 a constraint, ensuring alignment linear transformation preserves spectral clustering eigenvectors.
 70 Furthermore, the graph size is enormous as it is a product space over pixels, layers, images, and
 71 models; therefore, computing eigenvectors over their pairwise affinity matrix can be computationally
 72 infeasible. We developed a Nystrom-like approximation to ensure efficient computation.

73 In summary, our key contributions are:

- 74 1. We constructed a universal channel-aligned space using brain encoding as supervision and spectral clustering eigenvector constraints to ensure minimal channel signal loss. Brain encoding associates the aligned channel space to brain regions and gives them meanings.
- 75 2. Models trained with different objectives learned similar visual concepts: corresponding channel patterns exist across different models. The resulting visual concepts can be validated by unsupervised segmentation benchmarks on ImageNet-segmentation and PASCAL VOC.
- 76 3. Models show divergent computation paths over the visual concept space formed by the top-k spectral eigenvectors. Different models differ in trajectories and pace of movement layer-to-layer.

82 2 Methods: AlignedCut

83 Just as human languages might consist of distinct alphabets, features across different models appear superficially in embedding spaces as almost mutually orthogonal (Figure 3). However, the underlying information that they represent can be similar. To jointly analyze features across models and layers, we proposed the **channel align transform** that linearly projects features to a universal space.

91 The learning signal for the channel align transform is provided by **brain response prediction**. Learning from brain prediction offers two advantages.

94 First, brain response covers rich representations from all levels of semantics; the channel alignment removes irrelevant information while preserving the necessary and sufficient visual image features. 95 Second, knowledge of brain regions provides an interpretable understanding of their corresponding channels derived from the alignment.

98 Our visual concept discovery is formulated as a graph partitioning task using **spectral clustering**. We term our approach for this channel align and graph partitioning as **AlignedCut**. Furthermore, a major challenge in applying spectral clustering to large graphs is the complexity scaling issue. To address this, we developed a **Nystrom-like approximation** to reduce the computational complexity.

102 2.1 Brain-Guided Universal Channel Align

103 **Brain Dataset** We used the Algonauts competition (Gifford et al., 2023) release of Nature Scenes Dataset (NSD) (Allen et al., 2022). Briefly, NSD provides an fMRI brain scan when watching COCO images. Each subject viewed 10,000 images over 40 hours of scanning. We used the first subject’s publicly shared pre-processed and denoised (Prince et al., 2022) data.

107 **Channel Align** Let $\mathcal{V} = \{\mathbf{V}_1, \mathbf{V}_2, \dots, \mathbf{V}_n | \mathbf{V}_i \in \mathbb{R}^{P \times D_i}\}$ be the set of image features, extracted from each layer of pre-trained ViT models, where $P = (H \times W + 1)$ is image patches and class token, D_i is the hidden dimension. In particular, we used the attention layer output for each \mathbf{V}_i without adding residual connections from previous layers. Let \mathcal{V}' be the channel-aligned features; the goal of channel alignment is to learn a set of linear transform $\mathcal{W} = \{\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_n | \mathbf{W}_i \in \mathbb{R}^{D_i \times D'}\}$. In the new D' dimensional space, channels are aligned.

$$\mathcal{V}' = \mathcal{V} \odot \mathcal{W} = \{\mathbf{V}_1 \mathbf{W}_1, \mathbf{V}_2 \mathbf{W}_2, \dots, \mathbf{V}_n \mathbf{W}_n | \mathbf{V}_i \mathbf{W}_i \in \mathbb{R}^{P \times D'}\} \quad (1)$$

113 **Brain Prediction** To produce a learning signal for channel align \mathcal{W} , features from \mathcal{V}' are summed (not concatenated) to do brain prediction. Let $\mathbf{Y} \in \mathbb{R}^{1 \times N}$ be the brain prediction target, where N is the number of flattened 3D brain voxels, and 1 indicates that each voxel’s response is a scalar value. Let $F_\theta : \mathbb{R}^{P \times D'} \Rightarrow \mathbb{R}^{1 \times N}$ be the learned brain encoding model; without loss of generalizability, we set F_θ as global average pooling then linear weight $\beta_\theta \in \mathbb{R}^{D' \times N}$ and bias $\epsilon_\theta \in \mathbb{R}^{1 \times N}$:

$$\left[\text{AvgPool}_{p \in P} \left(\frac{1}{n} \sum_{i=1}^n (\mathbf{V}_i \mathbf{W}_i) \right) \times \beta_\theta + \epsilon_\theta \right] \Rightarrow \mathbf{Y} \quad (2)$$

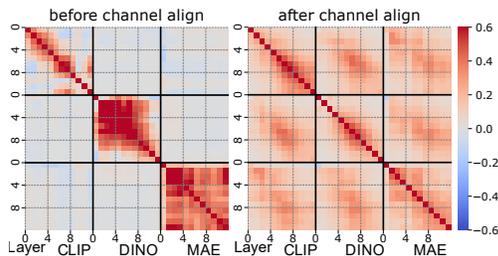


Figure 3: Cosine similarity of channel activation on the same image inputs.

118 **Channel in the Brain’s Space** Let $\mathcal{B} = \{B_1, B_2, \dots, B_n | B_i \in \mathbb{R}^{P \times N}\}$ be the set of channel
 119 activations in the brain’s space. By defining $B_i := V_i W_i \times \beta_\theta$, we have the brain response prediction
 120 $Y = \text{Avg Pool}_{p \in P} (\frac{1}{n} \sum_{i=1}^n B_i) + \epsilon_\theta$ (Eq. (2)). Intuitively, we linearly transformed the activation
 121 to the brain’s space, such that the activation from all slots sum up to the brain response prediction.

122 2.2 Graph Spectral Clustering

123 **Spectral Clustering** We use spectral clustering for visual concepts discovery and image-channel
 124 analysis; it provides 1) soft-cluster embedding space and 2) unsupervised hierarchical image seg-
 125 mentation. Normalized Cut (Shi and Malik, 2000) partitions the graph into sub-graphs with minimal
 126 cost of breaking edges. It embeds the graph into a lower dimensional eigenvector representation,
 127 where each eigenvector is a hierarchical sub-graph assignment.

128 Let $A \in \mathbb{R}^{M \times M}$ be the symmetric affinity matrix, where M denotes the total number of image
 129 patches. Given channel aligned features $V' \in \mathbb{R}^{M \times D'}$, we define $A_{ij} := \exp(\cos(V'_i, V'_j) - 1)$
 130 such that $A_{ij} > 0$ measures the similarity between data i and j . The spectral clustering embedding
 131 $X \in \mathbb{R}^{M \times C}$ is solved by the top C eigenvectors of the following generalized eigenproblem:

$$(D^{-1/2} A D^{-1/2}) X = X \Lambda \quad (3)$$

132 where D is the diagonal degree matrix $D_{ii} = \sum_j A_{ij}$, Λ is diagonal eigenvalue matrix.

133 **Nystrom-like Approximation** Computing eigenvectors for $A \in \mathbb{R}^{M \times M}$ is prohibitively expen-
 134 sive for enormous M with a time complexity of $O(M^3)$. The original Nystrom approximation
 135 method (Fowlkes et al., 2004) reduced the time complexity to $O(m^3 + m^2 M)$ by solving eigenvec-
 136 tors on sub-sampled graph $A' \in \mathbb{R}^{m \times m}$, where $m \ll M$. In particular, the orthogonalization step
 137 of eigenvectors introduced the time complexity of $O(m^2 M)$. Because our Nystrom-like approxi-
 138 mation trades the $O(m^2 M)$ orthogonalization term with the K-nearest neighbor, our Nystrom-like
 139 approximation reduced the time complexity to $O(m^3 + m M)$.

140 Our Nystrom-like Approximation first solves the eigenvector $X' \in \mathbb{R}^{m \times C}$ on a sub-sampled graph
 141 $A' \in \mathbb{R}^{m \times m}$ using Equation (3), then propagates the eigenvector from the sub-graph m nodes
 142 to the full-graph M nodes. Let $\tilde{X} \in \mathbb{R}^{M \times C}$ be the approximation $\tilde{X} \approx X$. The eigenvector
 143 approximation \tilde{X}_i of full-graph node $i \leq M$ is assigned by averaging the top K-nearest neighbors’
 144 eigenvector X'_k from the sub-graph nodes $k \leq m$:

$$\begin{aligned} \mathcal{K}_i &= KNN(A_{*i}; m, K) = \arg \max_{k \leq m} \sum_{k=1}^K A_{ki} \\ \tilde{X}_i &= \frac{1}{\sum_{k \in \mathcal{K}_i} A_{ki}} \sum_{k \in \mathcal{K}_i} A_{ki} X'_k \end{aligned} \quad (4)$$

145 where $KNN(A_{*i}; m, K)$ denotes KNN from full-graph node $i \leq M$ to sub-graph nodes $k \leq m$.

146 2.3 Affinity Eigen-constraints as Regularization for Channel Align

147 While brain prediction can provide strong supervi-
 148 sion for the learned channel align operation, we ob-
 149 served that the quality of unsupervised segmentation
 150 dropped after the channel alignment. To address this
 151 issue, a regularization term is added:

$$\mathcal{L}_{eigen} = \|X_b X_b^T - X_a X_a^T\| \quad (5)$$

152 where X_b and $X_a \in \mathbb{R}^{\tilde{m} \times c}$ are affinity matrix
 153 eigenvectors before and after channel alignment, re-
 154 spectively; $\tilde{m} = 100$ are randomly sampled nodes
 155 in a mini-batch and $c = 6$ are the top eigenvectors.
 156 The eigen-constraint preserves spectral clus-
 157 tering eigenvectors in dot-product space, invariant to random rotations in eigenvectors. We found
 158 adding eigen-constraints improved both the performance of segmentation (Figure 5) and the brain
 prediction score (Table 1).

Table 1: Affinity eigen-constraints improved brain score (R^2 : variance explained).

λ_{eigen}	ROI Brain Score R^2 (± 0.001)			
	V1	V4	EBA	all
1.0	0.170	0.181	0.295	0.196
0.1	0.167	0.179	0.294	0.193
0	0.155	0.166	0.296	0.188

159 **3 Results**

160 Our spectral clustering analysis aims to discover visual concepts that share the same pattern of
 161 channel activation across different models and layers. However, implementing spectral clustering
 162 analysis comes with two main challenges. First, the models sit in different feature spaces, so direct
 163 clustering will not reveal their overlap and similarities. Second, when scaling up to a large graph,
 164 spectral clustering is computationally expensive.

165 To address the first challenge, we developed our channel align transform to align features into a
 166 universal space. We extracted features from all 12 layers of the CLIP (ViT-B, OpenAI) (Radford
 167 et al., 2021), DINOv2 (ViT-B with registers) (Darcet et al., 2024), and MAE (ViT-B) (He et al.,
 168 2022) and then transformed features from each layer into the universal feature space.

169 To address the second challenge, we developed our Nystrom-like approximation to reduce the com-
 170 putational complexity. We extracted features from 1000 ImageNet (Deng et al., 2009) images, with
 171 each image consisting of 197 patches per layer. The entire product space of all images and fea-
 172 tures totaled $M = 7e+6$ nodes, from which we applied our Nystrom-like approximation with sub-
 173 sampled $m = 5e+4$ nodes and KNN $K = 100$, computing the top 20 eigenvectors.

174 To visualize the affinity eigenvectors, the top 20 eigenvectors were reduced to a 3-dimensional space
 175 by t-SNE, and a color value was assigned to each node by the RGB cube. We call this approach
 176 AlignedCut color.



Figure 4: Spectral clustering in the universal channel aligned feature space. The image pixels are colored by our approach AlignedCut, the pixel RGB value is assigned by the 3D spectral-tSNE of the top 20 eigenvectors. The coloring is consistent across all images, layers, and models.

- 177 In Figure 4, we displayed the analysis, AlignedCut color, and made the following observations:
- 178 1. In CLIP layer-5, DINO layer-6, and MAE layer-8, there is class-agnostic figure-ground separation, with foreground objects from different categories grouped into the same AlignedCut color.
 - 179 2. In CLIP layer-9, there is a class-specific separation of foreground objects, with foreground objects grouped into AlignedCut colors with associated semantic categories.
 - 180 3. Before layer-3, CLIP and DINO produce the same AlignedCut color regardless of the image
 - 181 182 183 input. From layer-4 onwards, the AlignedCut color smoothly changes over layers.

184 **3.1 Figure-ground representation emerge before categories**

185 In this section, we benchmark each layer in CLIP with unsupervised segmentation. The key findings
 186 from this benchmarking are: **1)** The figure-ground representation emerges at CLIP layer-4 and is
 187 preserved in subsequent layers; **2)** Categories emerge over layers, peaking at layer-9 and layer-10.

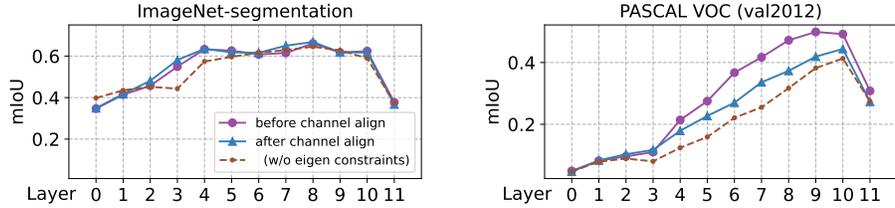


Figure 5: Unsupervised segmentation scores from spectral clustering on each CLIP layer. ImageNet-segmentation dataset is used with binary figure-ground labels, and the mIoU score peaks plateau from layer-4 to layer-10. In PASCAL VOC with 20 class labels, the mIoU score peaks at layer-9.

188 **From which layers did the figure-ground and category representations emerge?** We conducted
 189 experiments that compared the unsupervised segmentation scores across layers, tracing how well
 190 each representation is encoded at each layer. We used two datasets: a) ImageNet-segmentation
 191 (Guillaumin et al., 2014) with binary figure-ground labels, and b) PASCAL VOC (Everingham et al.,
 192 2010) with 20 category labels. The results are presented in Figure 5. On the ImageNet-segmentation
 193 benchmark, the score peaks at layer-4 (mIoU=0.6) and plateaus in subsequent layers, suggesting that
 194 the figure-ground representation is encoded and preserved from layer-4 onwards. On the PASCAL
 195 VOC benchmark, the score peaks at layer-9 and layer-10 (mIoU=0.5) even though it is low at layer-4
 196 (mIoU=0.2), indicating that category information is encoded at layer-9 and layer-10. Overall, we
 197 conclude that the figure-ground representation emerges before the category representation.

198 **3.2 Visual concepts: class-agnostic figure-ground**

199 In this section, we use brain activation heatmaps and image similarity heatmaps to describe figure-
 200 ground visual concepts. The key findings from these heatmaps are: **1)** The figure vs. ground pixels
 201 activate different channels; **2)** The figure-ground visual concept is class-agnostic; **3)** The figure-
 202 ground visual concept is consistent across models.

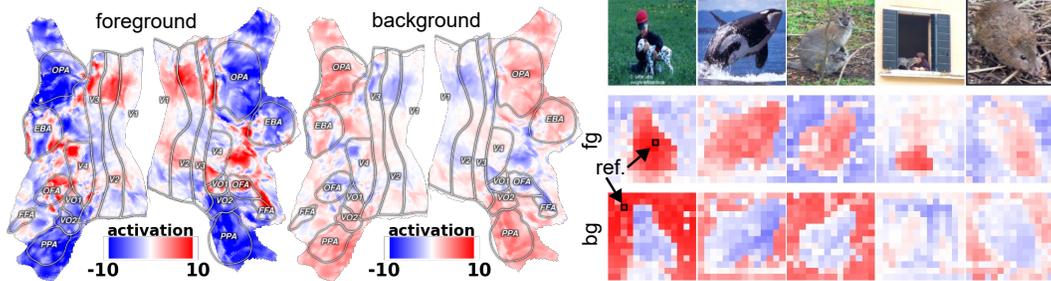


Figure 6: The figure-ground visual concepts in CLIP layer-5. **Left:** Mean activation of foreground or background pixels, linearly transformed to the brain’s space. **Right:** Cosine similarity from *one* reference pixel marked. The figure-ground visual concepts are agnostic to image categories.

203 **How can the channel activation patterns of the figure-ground visual concept be described?** We
 204 averaged the channel activations from foreground and background pixels, using the ground-truth
 205 labels from the ImageNet-segmentation dataset. The averaged channel activations were transformed
 206 into the brain’s space. In Figure 6, foreground pixels exhibit positive activation in early visual brain
 207 ROIs (V1 to V4) and the face-selective ROI (FFA), while negatively activating place-selective ROIs
 208 (OPA and PPA). Interestingly, background pixels activate the reverse pattern compared to foreground
 209 pixels. Overall, the figure and ground pixels activate distinct brain ROIs.

210 **Is the figure-ground visual concept class-agnostic?** We manually selected *one* pixel and computed
 211 the cosine similarity to all of the other image pixels. In Figure 6, the results demonstrate that one
 212 pixel (on the human) could segment out foreground objects from all other classes (shark, dog, cat,
 213 rabbit). The same result holds true for one background pixel. We conclude that the figure-ground
 214 visual concept is class-agnostic.

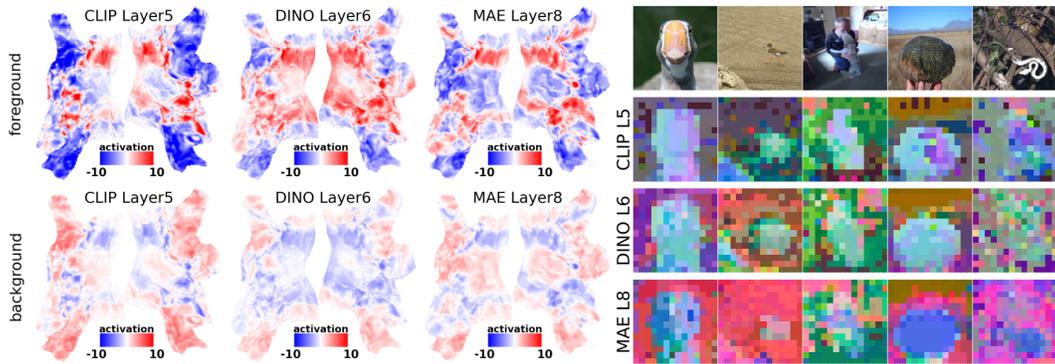


Figure 7: The same figure-ground visual concepts are found in CLIP, DINO and MAE. **Left:** Mean activation of all foreground (top) and background (bottom) pixels; the three models exhibit similar activation patterns. **Right:** AlignedCut, pixels colored by 3D spectral-tSNE of the top 20 eigenvectors; the three models show similar grouping colors for foreground pixels.

215 **Is the figure-ground visual concept consistent across models?** We performed the channel analysis
 216 for CLIP, DINO, and MAE. In Figure 7, the foreground or background pixels activates similar
 217 brain ROIs across the three models. Additionally, spectral clustering grouped the representations of
 218 foreground objects into similar colors for CLIP and DINO (light blue), the grouping for MAE is less
 219 similar (dark blue). Overall, the figure-ground visual concept is consistent across models.

220 3.3 Visual concepts: categories

221 In this section we use AlignedCut to discover category visual concepts. The key findings from the
 222 category visual concepts are: 1) Class-specific visual concepts activate diverse brain regions; 2)
 223 Visual concepts with higher channel activation values are more consistent.

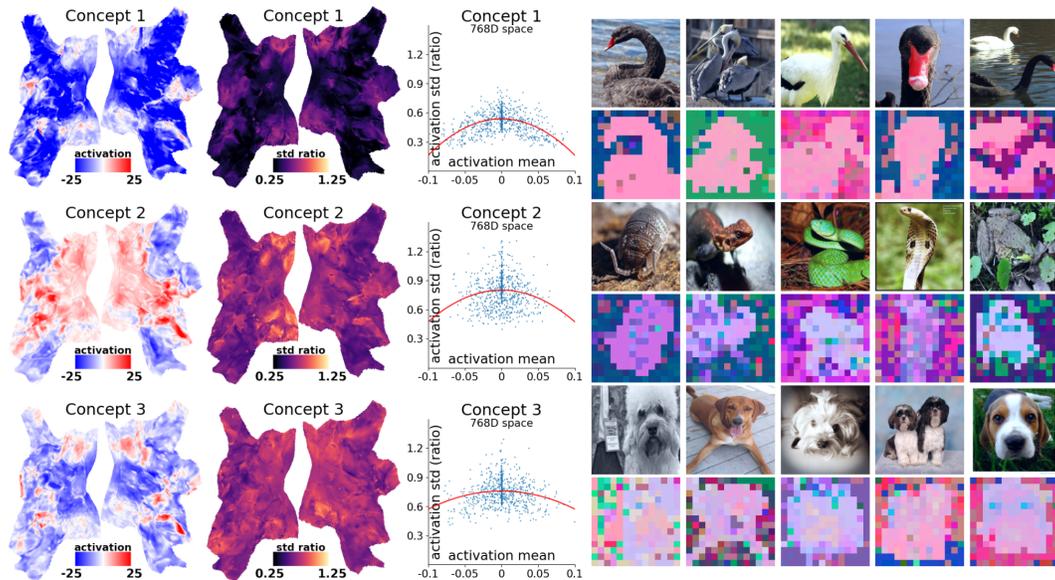


Figure 8: Category visual concepts in CLIP layer-9. **Left:** Mean activation of all pixels within an Euclidean sphere centered at the visual concept in the 3D spectral-tSNE space; the concepts activate different brain regions. **Middle:** The standard deviation negatively correlates with absolute mean activations. **Right:** AlignedCut, pixels colored by 3D spectral-tSNE of the top 20 eigenvectors.

224 **How does each class-specific concept activate the channels?** To answer this question, we sam-
 225 pled class-specific concepts from CLIP layer-9. First, we used farthest point sampling to identify
 226 candidate centers in the 3D spectral-tSNE space. Then, each candidate center was grouped with its
 227 neighboring pixels within an Euclidean sphere in the spectral-tSNE space. Finally, the channel acti-
 228 vations of the grouped pixels were averaged to produce the mean channel activation for each visual
 229 concept. In Figure 8, Concept 1 (duck, goose) negatively activates late brain regions; Concept 2

230 (snake, turtle) positively activates early brain regions and also FFA; Concept 3 (dog) negatively activates early brain regions. Overall, category-specific visual concepts activate diverse brain regions.
 231

232 **How do we quantify the consistency of each visual concept?** Qualitatively, Concept 1 exhibits more consistent coloring (Figure 8, pink) than Concept 3 (purple). To further quantify this observation, we
 233 computed the mean and standard deviation of channel activations for each Euclidean sphere centered on a concept. In Figure 8, there is a reverse U-shape relation between magnitude and standard deviation.
 234 The reverse U-shape implies that larger absolute mean channel activation corresponds to lower standard deviation. Overall, higher channel activation magnitudes suggest more consistent visual concepts.
 235
 236
 237
 238

239 3.4 Transition of visual concepts over layers

240 In this section, instead of using 3D spectral-tSNE, we use 2D spectral-tSNE to trace the layer-to-layer feature computation. The key findings of spectral-tSNE in 2D are: **1)** The figure vs. ground pixels are encoded in separate spaces in late layers; **2)** The representations for foreground and background bifurcate at CLIP layer-4 and DINO layer-5.
 241
 242
 243

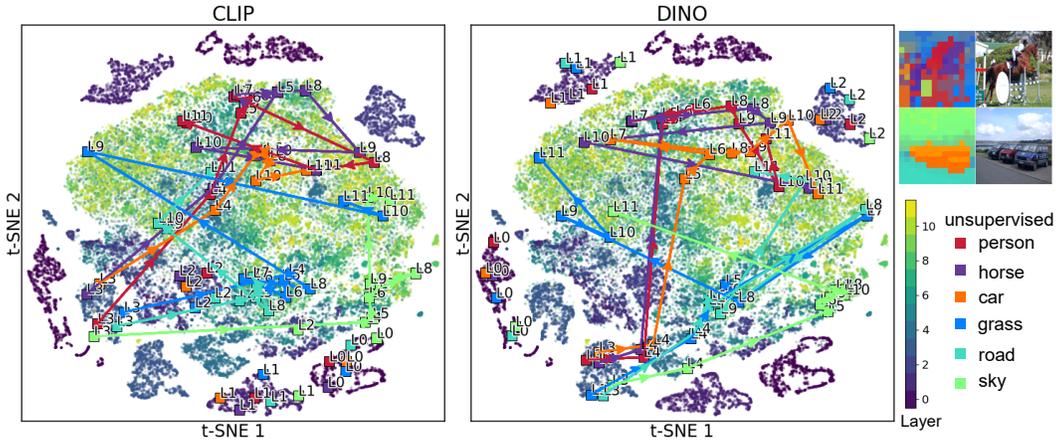


Figure 9: Trajectory of feature progression in layers for six example pixels. **Left:** 2D spectral-tSNE plot of the top 20 eigenvectors, jointly clustered across all models; the foreground and background pixels bifurcate at CLIP layer-4 and DINO layer-5. **Right:** Pixels colored by unsupervised segmentation.

244 **How does the network encode figure and ground pixels in each layer?** We performed spectral
 245 clustering and 2D t-SNE on the top 20 eigenvectors to project all layers into a 2D spectral-tSNE
 246 space. In Figure 9, we found that all foreground and background pixels are grouped together in
 247 each early layer. Each early layer (dark dots) forms an isolated cluster separate from other layers,
 248 while late layers (bright dots) are grouped in the center. In the late layers, there is a separation
 249 where foreground pixels occupy the upper part of 2D spectral-tSNE space, while background pixels
 250 occupy the middle part. Overall, foreground and background pixels are encoded in separate spaces
 251 in late layers.

252 **How does the network process each pixel from layer to layer?** In the 2D spectral-tSNE plot,
 253 we traced the trajectory for each pixel from layer-3 to the last layer. In Figure 9, we found that
 254 the trajectories for foreground and background pixels bifurcate: foreground pixels (person, horse,
 255 car) traverse to the upper side and remain within the upper side; background pixels (grass, road,
 256 sky) jump between the middle right and left sides. The same bifurcation is consistently observed
 257 for CLIP from layer-3 to layer-4 and DINO from layer-4 to layer-5. Furthermore, to quantify the
 258 bifurcation for foreground and background pixels, we first sampled 5 visual concepts from CLIP
 259 layer-3 and layer-4. Then, we measured the transition probability between visual concepts, defined
 260 as the proportion of pixels that transitioned from an Euclidean circle around concept A to a circle
 261 around concept B. In Figure 10, the transition probability of foreground pixels to the upper side
 262 (A1 to B0) is higher than that of background pixels (0.44 vs. 0.16), while the transition probability
 263 of background pixels to the right side (A4 to B4) is higher than that of foreground pixels (0.36 vs.
 264 0.06). Overall, this suggests a bifurcation of figure and ground pixel representations at the middle
 265 layers of both CLIP and DINO.

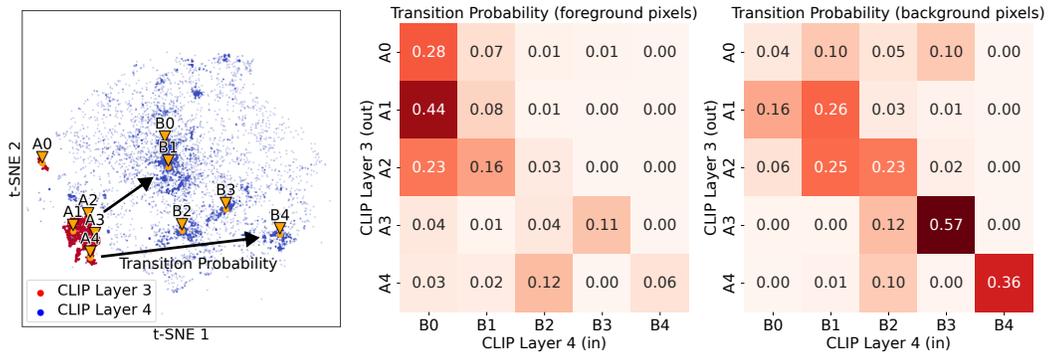


Figure 10: Transition probability of visual concepts from CLIP layer-3 to layer-4. **Left:** Five visual concepts sampled from CLIP layer-3 and layer-4. **Right:** Transition probability measured separately for foreground and background pixels; a bifurcation occurs where foreground pixels have more traffic to concept B0, while background pixels have more traffic to concepts B3 and B4.

266 4 Related Work

267 **Mechanistic Interpretability** is a field of study that intends to understand and explain the inner
 268 working mechanisms of deep networks. One approach is to interpret individual neurons (Bau et al.,
 269 2017; Dravid et al., 2023) and circuit connections between neurons (Olah et al., 2020). Another ap-
 270 proach is to interpret transformer attention heads (Gandelsman et al., 2024) and circuit connections
 271 between attention heads (Wang et al., 2023a). Other approaches also looked into the role of patch
 272 tokens (Sun et al., 2024). These approaches made the assumption that channels are aligned within
 273 the same model; we compare across models by actively aligning the channels to a universal space.

274 **Spectral Clustering** is a graphical method to analyze data grouping in the eigenvector space. Spec-
 275 tral methods have been widely used for unsupervised image segmentation (Shi and Malik, 2000;
 276 von Luxburg, 2007; Wu et al., 2018; Wang et al., 2023b). One major challenge for applying spectral
 277 clustering to large graphs is the complexity scaling issue. To solve the scaling issue, the Nystrom
 278 approximation (Fowlkes et al., 2004) approaches solve eigenvectors on sub-sampled graphs and then
 279 propagate to the full graph. Another approach is the gradient-based eigenvector solver (Zhang et al.,
 280 2023), which solves the eigenvectors in mini-batches. Our proposed Nystrom-like approximation
 281 achieves a computational speedup over the original Nystrom approximation, albeit at the expense of
 282 weakened orthogonality of the eigenvectors.

283 **Brain Encoding Model** is widely used by the computational neuroscience community (Kriegeskorte
 284 and Douglas, 2018). They have been using deep nets to explain the brain’s function. One approach
 285 is to use the gradient of the brain encoding model to find the most salient image features (Sarch
 286 et al., 2023). Another approach generate text caption for brain activation (Luo et al., 2024). Other
 287 approaches compare brain prediction performance for different models (Schrimpf et al., 2020). The
 288 field focused on using deep nets as a tool to explain the brain’s function; we go in the opposite
 289 direction by using the brain to explain deep nets.

290 5 Conclusion and Limitations

291 We present a novel approach to interpreting deep neural networks by leveraging brain data. Our
 292 fundamental innovation is twofold: First, we use brain prediction as guidance to align channels from
 293 different models into a universal feature space; Second, we developed a Nystrom-like approximation
 294 to scale up the spectral clustering analysis. Our key discovery is that recurring visual concepts exist
 295 across networks and layers; such concepts correspond to different levels of objects, ranging from
 296 figure-ground to categories. Additionally, we quantified the information flow from layer to layer,
 297 where we found a bifurcation of figure-ground visual concepts.

298 **Limitations.** While the learned channel align transformation projects all features onto a universal
 299 feature space, the nature of learned transformation does not preserve all the information. There
 300 is a small drop in unsupervised segmentation performance after channel alignment, which is not
 301 fully addressed by our proposed eigen-constraint regularization. Secondly, as a trade-off for faster
 302 computation, our Nystrom-like approximation does not produce strictly orthogonal eigenvectors. To
 303 produce expressive eigenvectors, our approximation relies on using larger sub-sample sizes than the
 304 original Nystrom method.

305 References

- 306 Allen, E. J., St-Yves, G., Wu, Y., Breedlove, J. L., Prince, J. S., Dowdle, L. T., Nau, M., Caron,
307 B., Pestilli, F., Charest, I., Hutchinson, J. B., Naselaris, T., and Kay, K. (2022). A massive 7T
308 fMRI dataset to bridge cognitive neuroscience and artificial intelligence. *Nature Neuroscience*,
309 25(1):116–126. Number: 1 Publisher: Nature Publishing Group. 3
- 310 Bau, D., Zhou, B., Khosla, A., Oliva, A., and Torralba, A. (2017). Network Dissection: Quantifying
311 Interpretability of Deep Visual Representations. In *Computer Vision and Pattern Recognition*. 9
- 312 Darcet, T., Oquab, M., Mairal, J., and Bojanowski, P. (2024). Vision Transformers Need Registers.
313 In *The Twelfth International Conference on Learning Representations*. 5
- 314 Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). ImageNet: A large-
315 scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern
316 Recognition*, pages 248–255. 5
- 317 Dravid, A., Gandselman, Y., Efros, A. A., and Shocher, A. (2023). Rosetta Neurons: Mining the
318 Common Units in a Model Zoo. In *Proceedings of the IEEE/CVF International Conference on
319 Computer Vision (ICCV)*, pages 1934–1943. 1, 9
- 320 Everingham, M., Van Gool, L., Williams, C. K. I., Winn, J., and Zisserman, A. (2010). The Pascal
321 Visual Object Classes (VOC) Challenge. *International Journal of Computer Vision*, 88(2):303–
322 338. 6
- 323 Fowlkes, C., Belongie, S., Chung, F., and Malik, J. (2004). Spectral grouping using the Nystrom
324 method. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(2):214–225. 4, 9
- 325 Gandselman, Y., Efros, A. A., and Steinhardt, J. (2024). Interpreting CLIP’s Image Representation
326 via Text-Based Decomposition. In *The Twelfth International Conference on Learning Represen-
327 tations*. 9
- 328 Gifford, A. T., Lahner, B., Saba-Sadiya, S., Vilas, M. G., Lascelles, A., Oliva, A., Kay, K., Roig, G.,
329 and Cichy, R. M. (2023). The Algonauts Project 2023 Challenge: How the Human Brain Makes
330 Sense of Natural Scenes. arXiv:2301.03198 [cs, q-bio]. 3
- 331 Guillaumin, M., Küttel, D., and Ferrari, V. (2014). ImageNet Auto-Annotation with Segmentation
332 Propagation. *International Journal of Computer Vision*, 110(3):328–348. 6
- 333 He, K., Chen, X., Xie, S., Li, Y., Dollár, P., and Girshick, R. (2022). Masked Autoencoders Are
334 Scalable Vision Learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
335 Pattern Recognition (CVPR)*, pages 16000–16009. 5
- 336 Kriegeskorte, N. and Douglas, P. K. (2018). Cognitive computational neuroscience. *Nature Neuro-
337 science*, 21(9):1148–1160. Publisher: Nature Publishing Group. 9
- 338 Lin, T.-Y., Goyal, P., Girshick, R., He, K., and Dollar, P. (2017). Focal Loss for Dense Object
339 Detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. 14
- 340 Luo, A., Henderson, M. M., Tarr, M. J., and Wehbe, L. (2024). BrainSCUBA: Fine-Grained Natural
341 Language Captions of Visual Cortex Selectivity. In *The Twelfth International Conference on
342 Learning Representations*. 9
- 343 Olah, C., Cammarata, N., Schubert, L., Goh, G., Petrov, M., and Carter, S. (2020). Zoom In: An
344 Introduction to Circuits. *Distill*, 5(3):e00024.001. 9
- 345 Prince, J. S., Charest, I., Kurzawski, J. W., Pyles, J. A., Tarr, M. J., and Kay, K. N. (2022). Im-
346 proving the accuracy of single-trial fMRI response estimates using GLMsingle. *eLife*, 11:e77599.
347 Publisher: eLife Sciences Publications, Ltd. 3
- 348 Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A.,
349 Mishkin, P., Clark, J., Krueger, G., and Sutskever, I. (2021). Learning Transferable Visual Models
350 From Natural Language Supervision. In Meila, M. and Zhang, T., editors, *Proceedings of the 38th
351 International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning
352 Research*, pages 8748–8763. PMLR. 5
- 353 Sarch, G. H., Tarr, M. J., Fragkiadaki, K., and Wehbe, L. (2023). Brain Dissection: fMRI-trained
354 Networks Reveal Spatial Selectivity in the Processing of Natural Images. In *Thirty-seventh Con-
355 ference on Neural Information Processing Systems*. 9

- 356 Schrimpf, M., Kubilius, J., Lee, M. J., Murty, N. A. R., Ajemian, R., and DiCarlo, J. J. (2020). Inte-
357 grative Benchmarking to Advance Neurally Mechanistic Models of Human Intelligence. *Neuron*.
358 9
- 359 Shi, J. and Malik, J. (2000). Normalized cuts and image segmentation. *IEEE Transactions on*
360 *Pattern Analysis and Machine Intelligence*, 22(8):888–905. 4, 9
- 361 Sun, M., Chen, X., Kolter, J. Z., and Liu, Z. (2024). Massive Activations in Large Language Models.
362 arXiv:2402.17762 [cs]. 9
- 363 von Luxburg, U. (2007). A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416.
364 9
- 365 Wang, K. R., Variengien, A., Conmy, A., Shlegeris, B., and Steinhardt, J. (2023a). Interpretabil-
366 ity in the Wild: a Circuit for Indirect Object Identification in GPT-2 Small. In *The Eleventh*
367 *International Conference on Learning Representations*. 9
- 368 Wang, X., Girdhar, R., Yu, S. X., and Misra, I. (2023b). Cut and learn for unsupervised object
369 detection and instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer*
370 *Vision and Pattern Recognition*, pages 3124–3134. 9
- 371 Wu, Z., Xiong, Y., Stella, X. Y., and Lin, D. (2018). Unsupervised Feature Learning via Non-
372 Parametric Instance Discrimination. In *Proceedings of the IEEE Conference on Computer Vision*
373 *and Pattern Recognition*. 9
- 374 Yang, H., Gee, J., and Shi, J. (2024). Brain Decodes Deep Nets. arXiv:2312.01280 [cs]. 1
- 375 Zhang, X., Yunis, D., and Maire, M. (2023). Deciphering ‘What’ and ‘Where’ Visual Pathways
376 from Spectral Clustering of Layer-Distributed Neural Representations. arXiv:2312.06716 [cs]. 9

377 **A Appendix overview**

- 378 1. Appendix **B** summarizes background of brain ROIs.
- 379 2. Appendix **C** is implementation details
- 380 2.1. Additional regularization terms
- 381 2.2. Brain encoding model training loss function
- 382 2.3. Unsupervised segmentation evaluation pipeline
- 383 2.4. Nystrom-like approximation for t-SNE
- 384 3. Appendix **D** lists more image examples from the 3D spectral-tSNE.
- 385 4. Appendix **E** lists figure-ground channel activation for every model and layer.
- 386 5. Appendix **F** lists more example category-specific visual concepts.
- 387 6. Appendix **G** lists more example pixels from the 2D spectral-tSNE information flow.

388 **B Brain Region Background Knowledge**

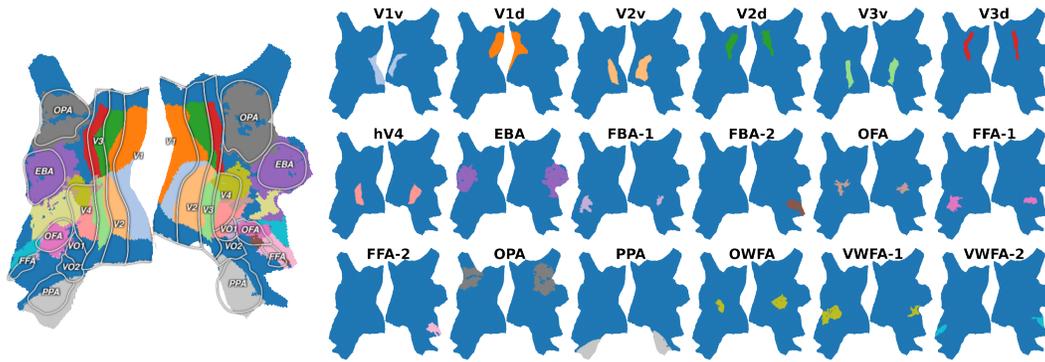


Figure 11: **Brain Region of Interests (ROIs)**. V1v: ventral stream, V1d: dorsal stream.

Table 2: Known function and selectivity of brain region of interests (ROIs).

ROI name	V1	V2	V3	V4	EBA	FBA	OFA	FFA	OPA	PPA	OWFA	VWFA
Known Function/Selectivity	primary visual			mid-level	body		face		navigation	scene	words	

389 This section briefly summarizes the known functions of key brain regions of interest (ROIs). Fig-
 390 ure 11 provides an overview of these brain ROIs. Table 2 lists the known functions and selectivities
 391 for each ROI.

392 In brief, V1 to V3 are the primary visual stream, which is further divided into ventral (lower) and dor-
 393 sal (upper) streams. V4 is a mid-level visual area. EBA (extrastriate body area) and FBA (fusiform
 394 body area) are selectively responsive to bodies, while FFA (fusiform face area) and OFA (occipital
 395 face area) show selectivity for faces. OWFA (occipital word form area) and VWFA (visual word
 396 form area) are selective for written words. PPA (parahippocampal place area) exhibits selectivity for
 397 scenes and places, and OPA (occipital place area) is involved in navigation and spatial reasoning.

398 Visual information processing in the brain follows a hierarchical, feedforward organization. Be-
 399 ginning in the primary visual cortex (V1) and progressing through higher visual areas like V2, V3,
 400 and V4, neurons exhibit increasingly large receptive fields and represent increasingly abstract visual
 401 concepts. While neurons in V1 encode low-level features like edges and orientations within a small
 402 portion of the visual field, neurons in V4 synthesize more complex patterns and object representa-
 403 tions across a larger area of the visual input.

404 C Implementation Details

405 C.1 Additional Regularization for Channel Align Transformation

406 Additional Regularization are added to the channel align transform to ensure good properties of the
407 aligned features: 1) zero-centered, 2) small covariance between channels, and 3) focal loss.

408 **Zero-centered regularization.** We did not apply z-score normalization to the extracted features;
409 instead, we added a regularization term to ensure the transformed features are zero-centered. Recall
410 that the channel-aligned transformed feature $\mathbf{V}' \in \mathbb{R}^{M \times D'}$, where M is the number of data points
411 and D' is the hidden dimension. The zero-center loss is defined as:

$$\mathcal{L}_{\text{zero}} = \frac{1}{D'} \frac{1}{M} \sum_{i \leq M, j \leq D'} v'_{ij} \quad (6)$$

412 **Covariance regularization.** We used the covariance loss to minimize the off-diagonal elements in
413 the covariance matrix of the transformed feature $C(\mathbf{V}')$, aiming to bring them close to $\mathbf{0}$. Recall
414 that channel align transformed feature $\mathbf{V}' \in \mathbb{R}^{M \times D'}$, where M is number of data, D' is the hidden
415 dimension. The covariance loss is defined as:

$$\mathcal{L}_{\text{cov}} = \frac{1}{D'} \sum_{i \neq j} [C(\mathbf{V}')]_{i,j}^2, \text{ where } C(\mathbf{V}') = \frac{1}{M-1} \sum_{i=1}^M (v'_i - \bar{v}') (v'_i - \bar{v}')^T, \bar{v}' = \frac{1}{M} \sum_{i=1}^M v'_i. \quad (7)$$

416 **Focal Loss.** Lin et al. (2017) introduced focal loss, which dynamically assigns smaller weights to
417 the loss function for hard-to-classify classes. In our scenario, we apply spectral clustering on the
418 affinity matrix $\mathbf{A}_a \in \mathbb{R}^{M \times M}$ after performing the channel alignment transform, where M represents
419 the number of data points. Due to the characteristics of spectral clustering, disconnected edges
420 play a more critical role than connected edges. Adding an edge between disconnected clusters
421 significantly reshapes the eigenvectors, while adding edges to connected clusters has only a minor
422 impact. Therefore, we aim to assign larger weights to disconnected edges in the loss function:

$$\mathcal{L}_{\text{eigen}} = \|(\mathbf{X}_b \mathbf{X}_b^T - \mathbf{X}_a \mathbf{X}_a^T) * \exp(-\mathbf{A}_b)\| \quad (8)$$

423 where $\mathbf{A}_b \in \mathbb{R}^{M \times M}$ is the affinity matrix before the channel alignment transform, element wise dot-
424 product to $\exp(-\mathbf{A}_b)$ assigned larger weights for disconnected edges. $\mathbf{X}_b \in \mathbb{R}^{M \times C}$, $\mathbf{X}_a \in \mathbb{R}^{M \times C}$
425 are eigenvectors before and after channel align transform, respectively.

426 C.2 Brain Encoding Model Training Loss

427 Let $\mathbf{Y} \in \mathbb{R}^{1 \times N}$ represent the brain prediction target, where N is the number of flattened 3D brain
428 voxels, and the 1 indicates that each voxel’s response is a scalar value. $\hat{\mathbf{Y}}$ is the model’s predicted
429 brain response. The brain encoding model training loss is the L1 loss:

$$\mathcal{L}_{\text{brain}} = \|\mathbf{Y} - \hat{\mathbf{Y}}\| \quad (9)$$

430 C.3 Total Training Loss

431 The total training loss is a combination of the following components: 1) brain encoding model loss,
432 2) eigen-constraint regularization, 3) zero-centered regularization, and 4) covariance regularization:

$$\mathcal{L} = \mathcal{L}_{\text{brain}} + \lambda_{\text{eigen}} \mathcal{L}_{\text{eigen}} + \lambda_{\text{zero}} \mathcal{L}_{\text{zero}} + \lambda_{\text{cov}} \mathcal{L}_{\text{cov}} \quad (10)$$

433 where we set $\lambda_{\text{eigen}} = 1$, $\lambda_{\text{zero}} = 0.01$, $\lambda_{\text{cov}} = 0.01$.

434 **C.4 Oracle-based Unsupervised Segmentation Evaluation Pipeline**

435 Our unsupervised segmentation pipeline aims to benchmark and compare the performance across
436 each single layer of the CLIP model. The evaluation pipeline is oracle-based:

- 437 1. Apply spectral clustering jointly across all images, taking the top 10 eigenvectors.
- 438 2. For each class of object (plus one background class), use ground-truth labels from the dataset
439 to mask out the pixels and their eigenvectors, and then use the mean of the eigenvectors to define a
440 center for each class.
- 441 3. Compute the cosine similarity of each pixel to all class centers.
- 442 4. For each pixel, if the maximum similarity to all classes is less than a threshold value, assign
443 this pixel to the background class.
- 444 5. Assign pixels (with a similarity greater than the threshold value) to the class with the maximum
445 similarity.

446 There’s one hyper-parameter, the threshold value that requires different optimal value for each layer
447 of CLIP. To ensure a fair comparison across all layers, the threshold value is grid-searched from 10
448 evenly spaced values between 0 and 1, the maximum mIoU score in the grid search is taken for each
449 layer.

450 **C.5 Nystrom-like approximation for t-SNE**

451 To visualize the eigenvectors, we applied t-SNE to the eigenvectors $\mathbf{X} \in \mathbb{R}^{M \times C}$, where the number
452 of data points M span the product space of models, layers, pixels, and images. Due to the enormous
453 size of $M = 7e+6$ nodes, t-SNE suffered from complexity scaling issues. We again applied our
454 Nystrom-like approximation to t-SNE, with sub-sampled $m = 10e+4$ nodes and KNN $K = 1$.

455 It’s worth noting that, since the non-linear distance adjustment in t-SNE, it’s crucial to use only one
456 nearest neighbor $K = 1$ for t-SNE.

457 **C.6 Computation Resource**

458 All of our experiments are performed on one consumer-grade RTX 4090 GPU. The brain encoding
459 model training took 3 hours on 4GB of VRAM, spectral clustering eigen-decomposition on large
460 graph took 10 minutes on 10GB of VRAM and 60GB of CPU RAM.

461 **C.7 Code Release**

462 Our code will be publicly released upon publication.

463 D 3D spectral-tSNE

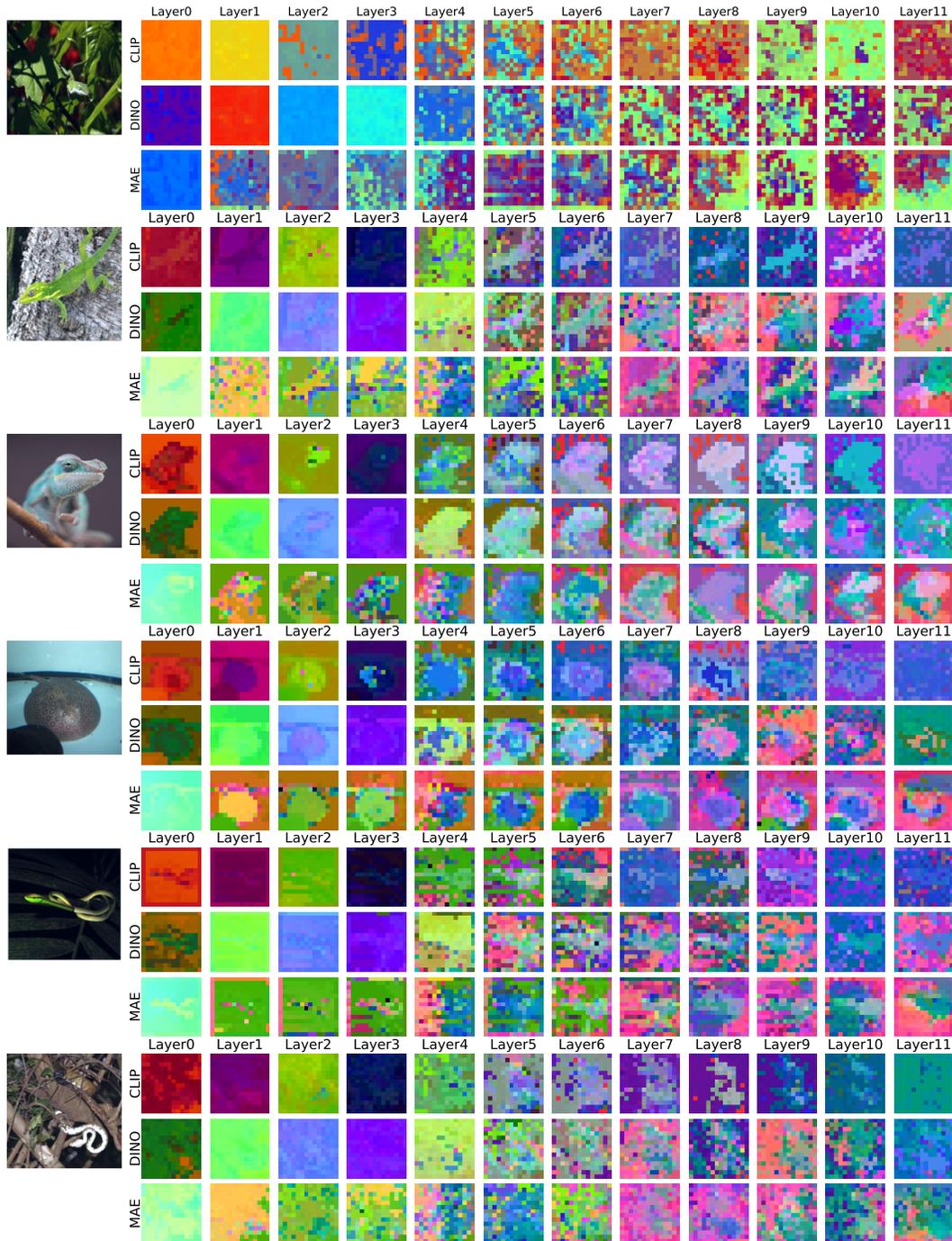


Figure 12: Spectral clustering in the universal channel aligned feature space. The image pixels are colored by our approach AlignedCut, the pixel RGB value is assigned by the 3D spectral-tSNE of the top 20 eigenvectors. The coloring is consistent across all images, layers, and models.

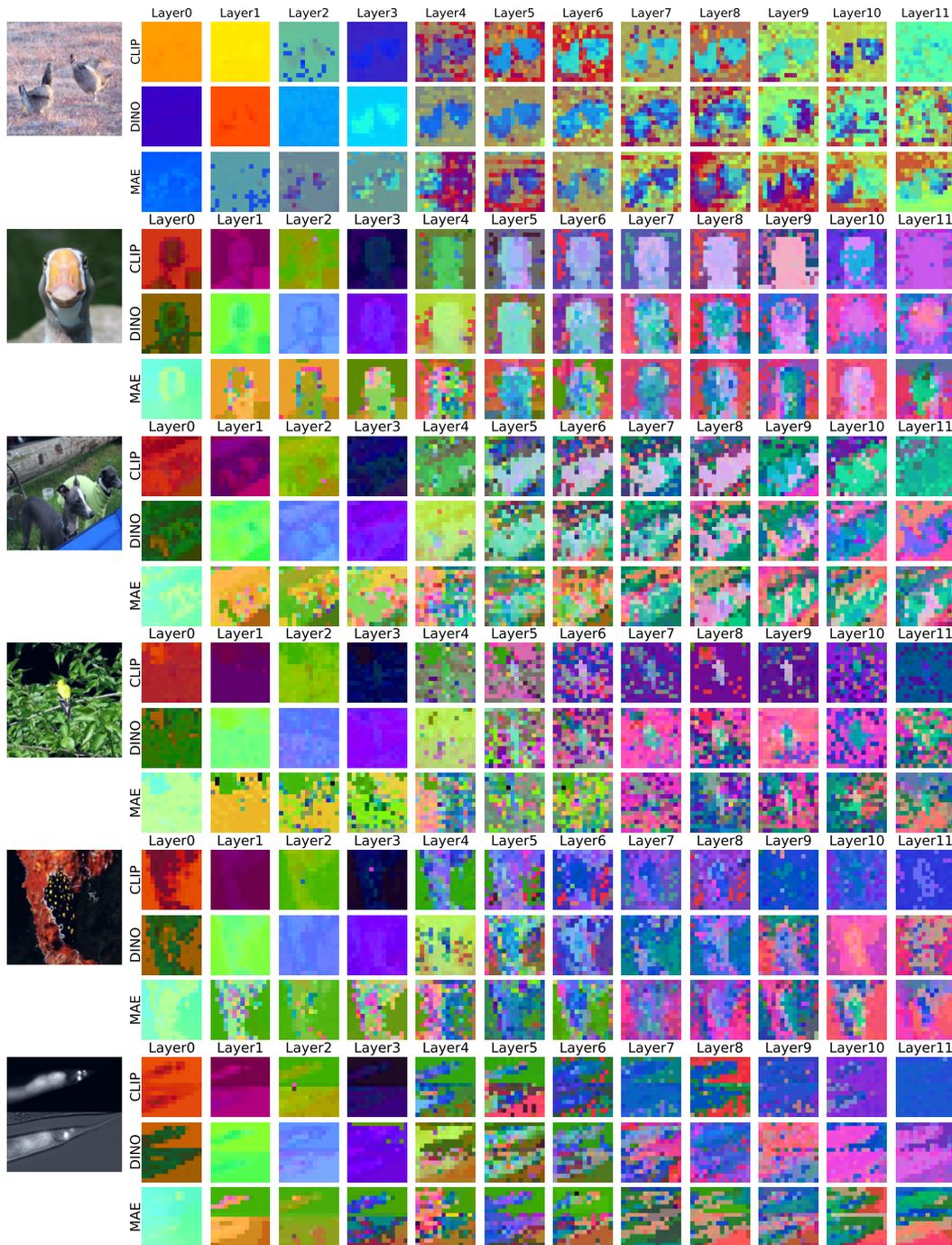


Figure 13: Spectral clustering in the universal channel aligned feature space. The image pixels are colored by our approach AlignedCut, the pixel RGB value is assigned by the 3D spectral-tSNE of the top 20 eigenvectors. The coloring is consistent across all images, layers, and models.

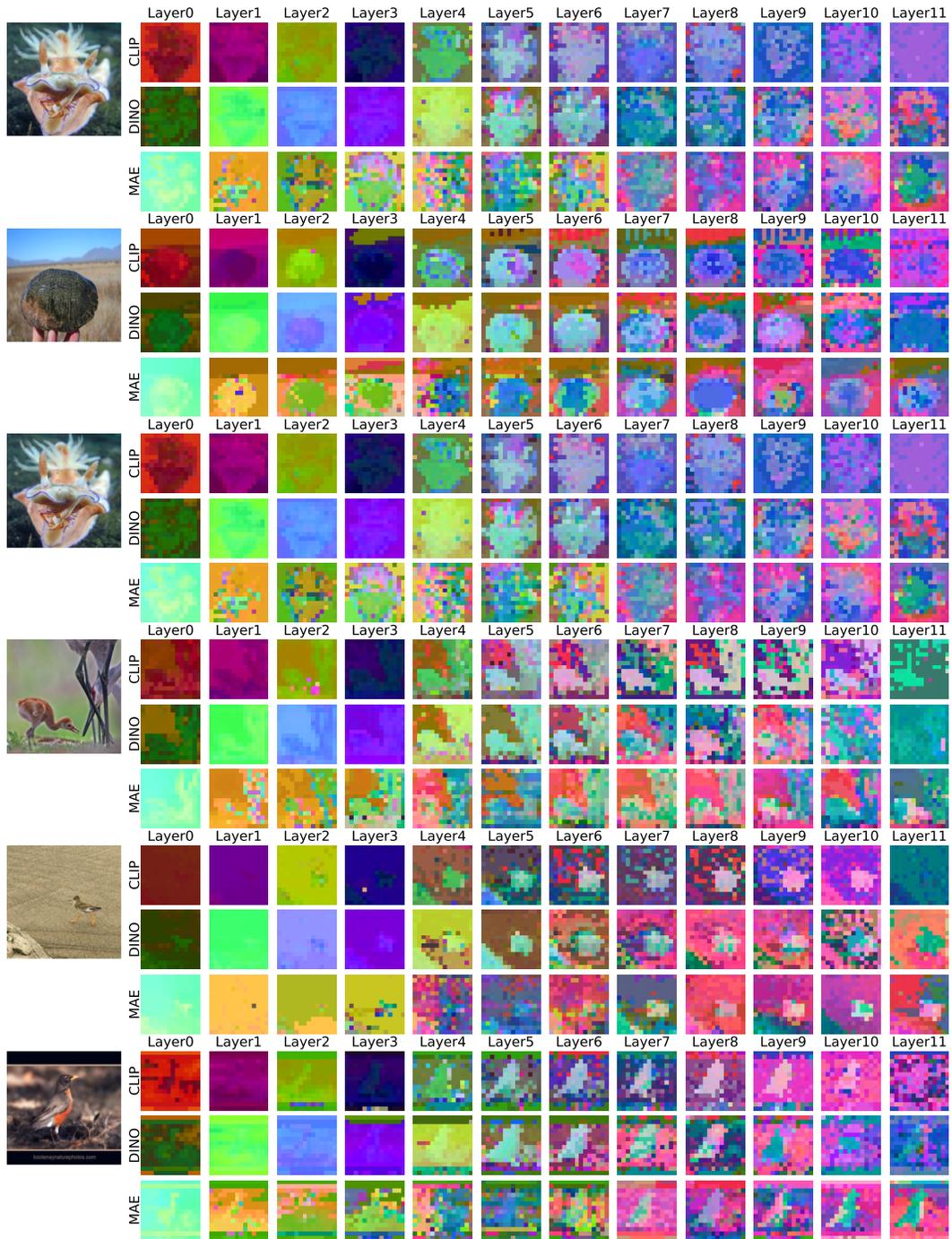


Figure 14: Spectral clustering in the universal channel aligned feature space. The image pixels are colored by our approach AlignedCut, the pixel RGB value is assigned by the 3D spectral-tSNE of the top 20 eigenvectors. The coloring is consistent across all images, layers, and models.

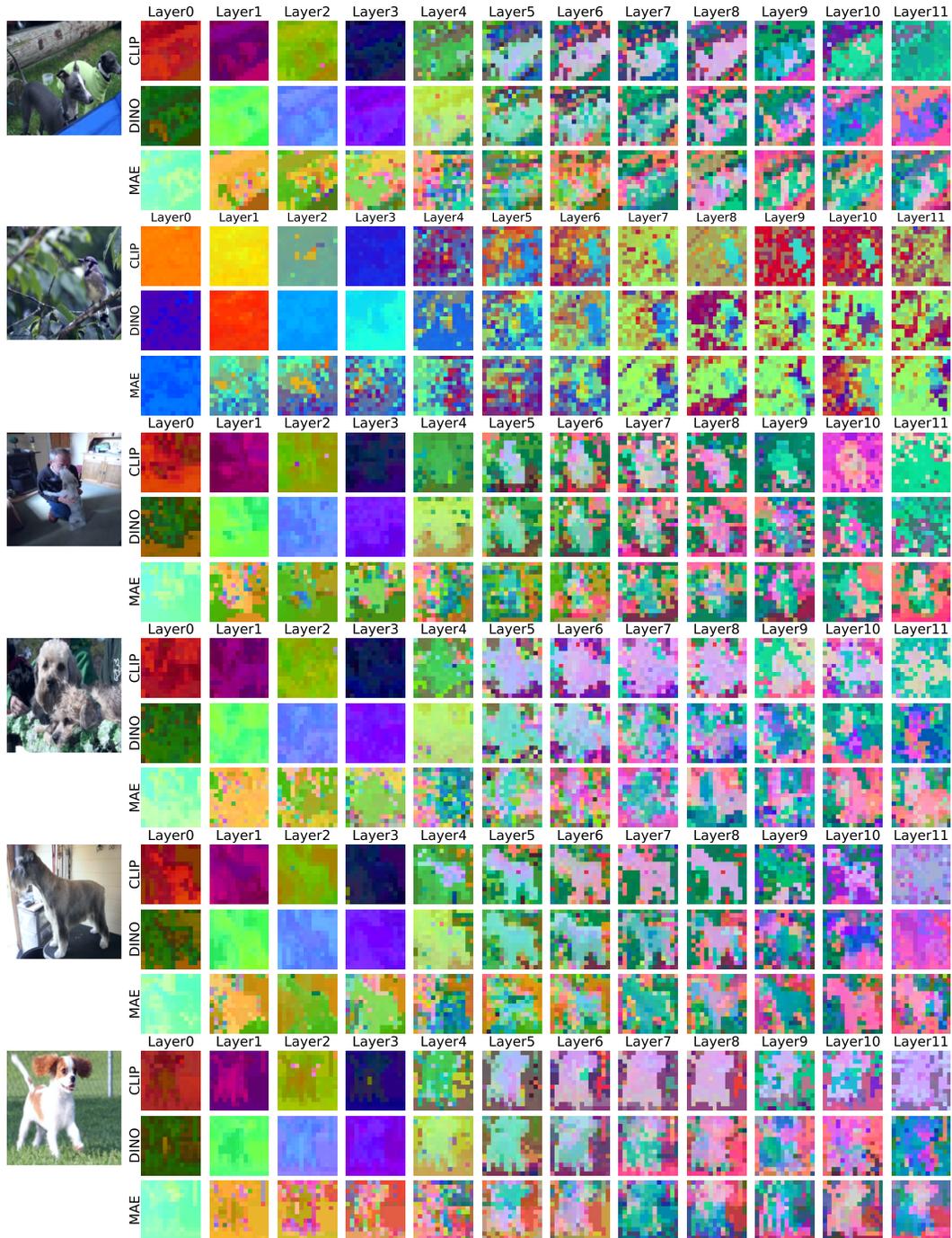


Figure 15: Spectral clustering in the universal channel aligned feature space. The image pixels are colored by our approach AlignedCut, the pixel RGB value is assigned by the 3D spectral-tSNE of the top 20 eigenvectors. The coloring is consistent across all images, layers, and models.

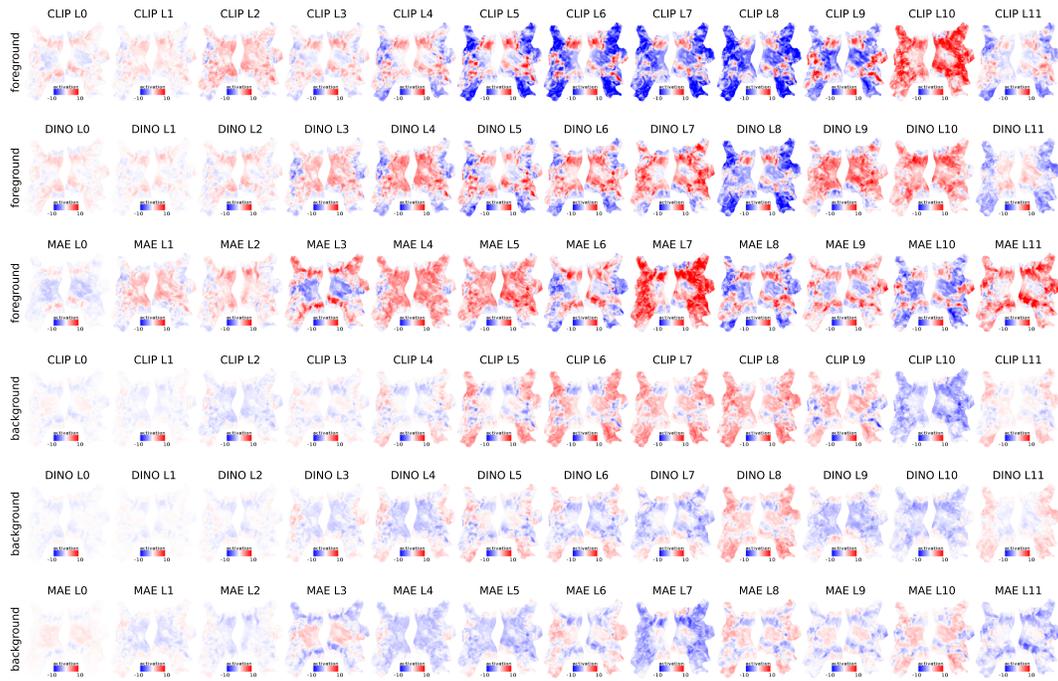
464 **E Figure-ground Channel Activation from All Layers and Models**

Figure 16: Mean activation of foreground or background pixels at each layer of CLIP, DINO and MAE. Channel activations are linearly transformed to the brain's space. Large absolute activation value means more consistent visual concepts.

465 **F Visual Concepts: Categories**

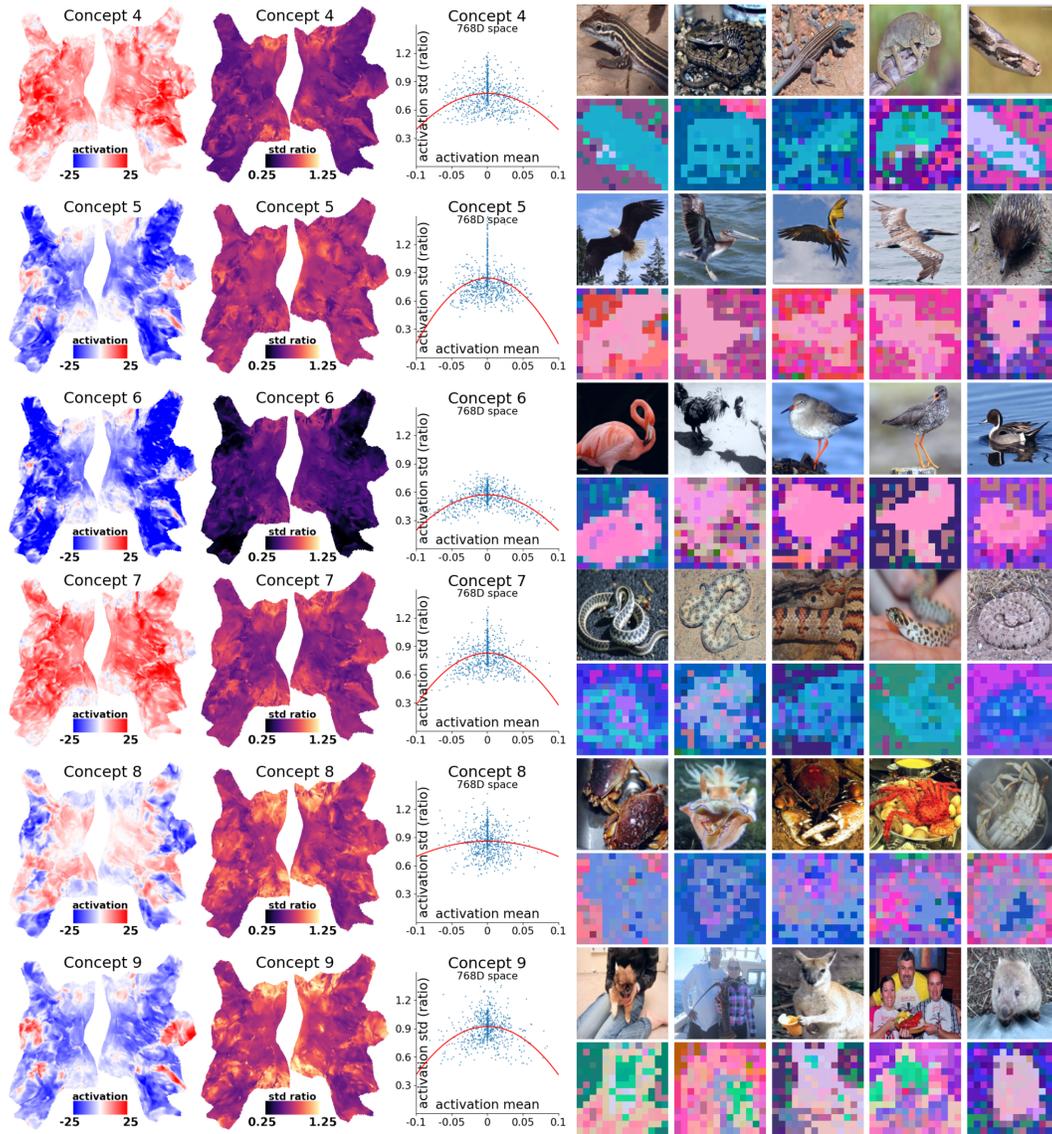


Figure 17: Category visual concepts in CLIP Layer 9. **Left:** Mean activation of all pixels within an Euclidean sphere centered at the visual concept in the 3D spectral-tSNE space; the concepts activate different brain regions. **Middle:** The standard deviation negatively correlates with absolute mean activations. **Right:** Spectral clustering, colored by 3D spectral-tSNE of the top 20 eigenvectors.

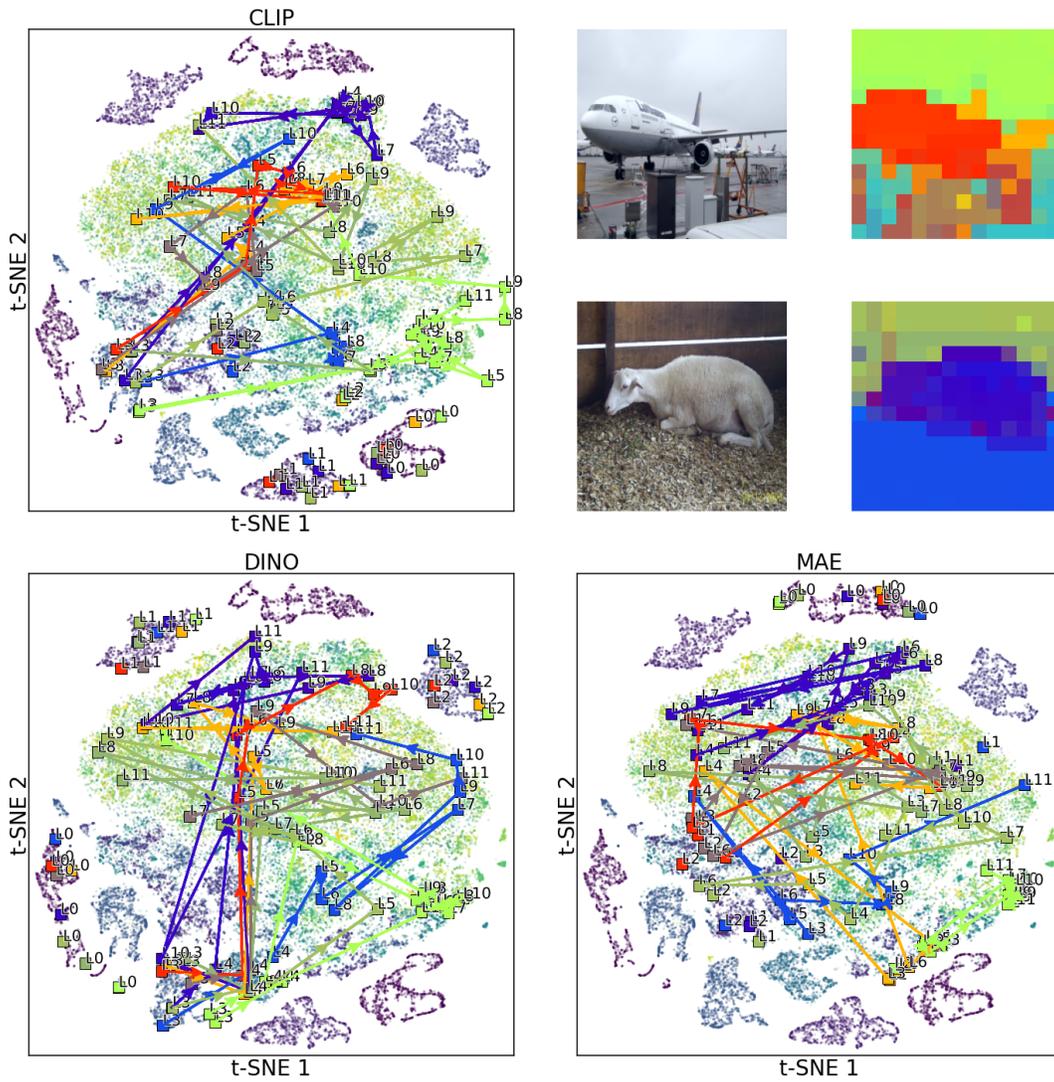
466 **G Layer-to-Layer Feature Computation Flow in 2D spectral-tSNE space**

Figure 18: Trajectory of feature progression in from layer to layer, in the 2D spectral-tSNE space. Arrows displayed for 10 randomly sampled example pixels. **Top Right:** Pixels are colored by unsupervised segmentation.

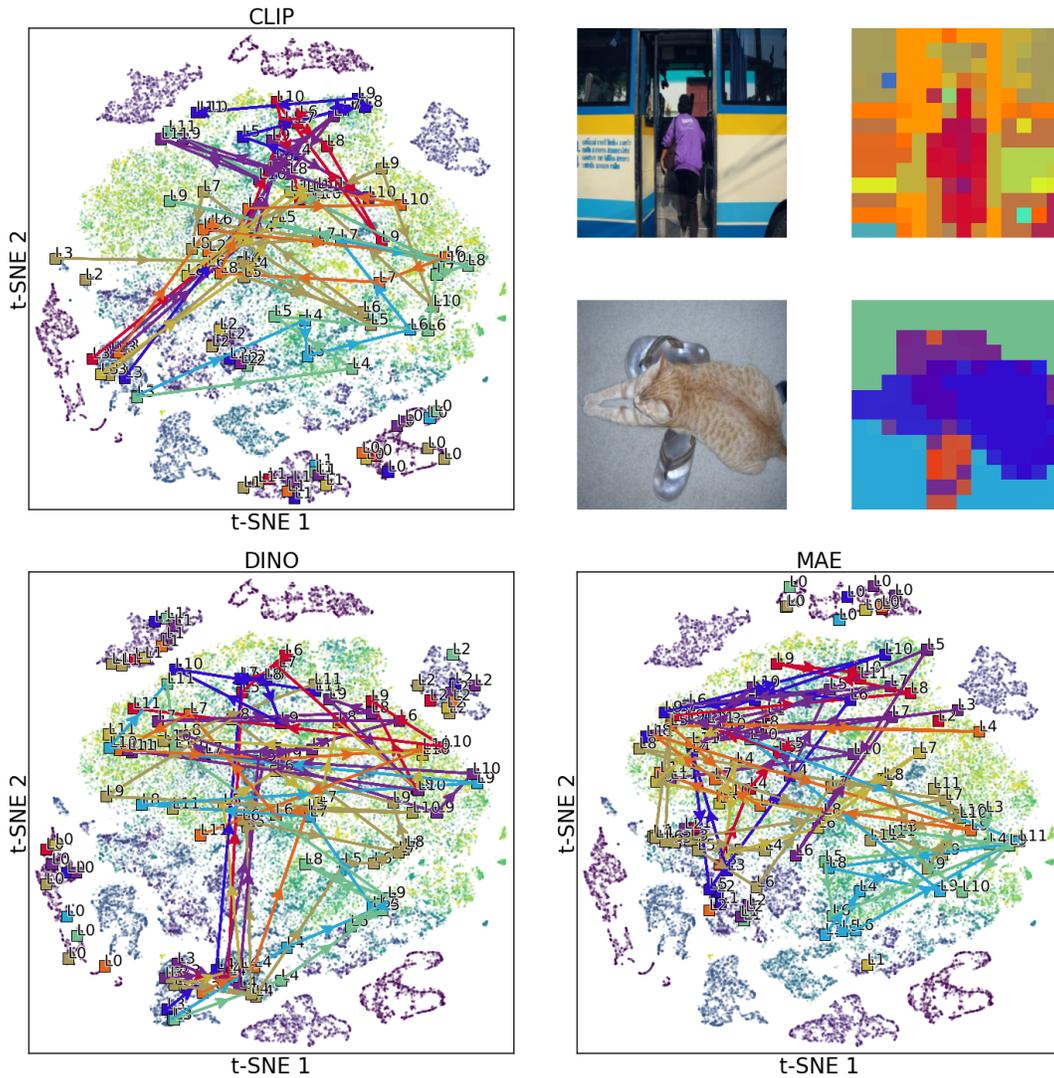


Figure 19: Trajectory of feature progression in from layer to layer, in the 2D spectral-tSNE space. Arrows displayed for 10 randomly sampled example pixels. **Top Right:** Pixels are colored by unsupervised segmentation.

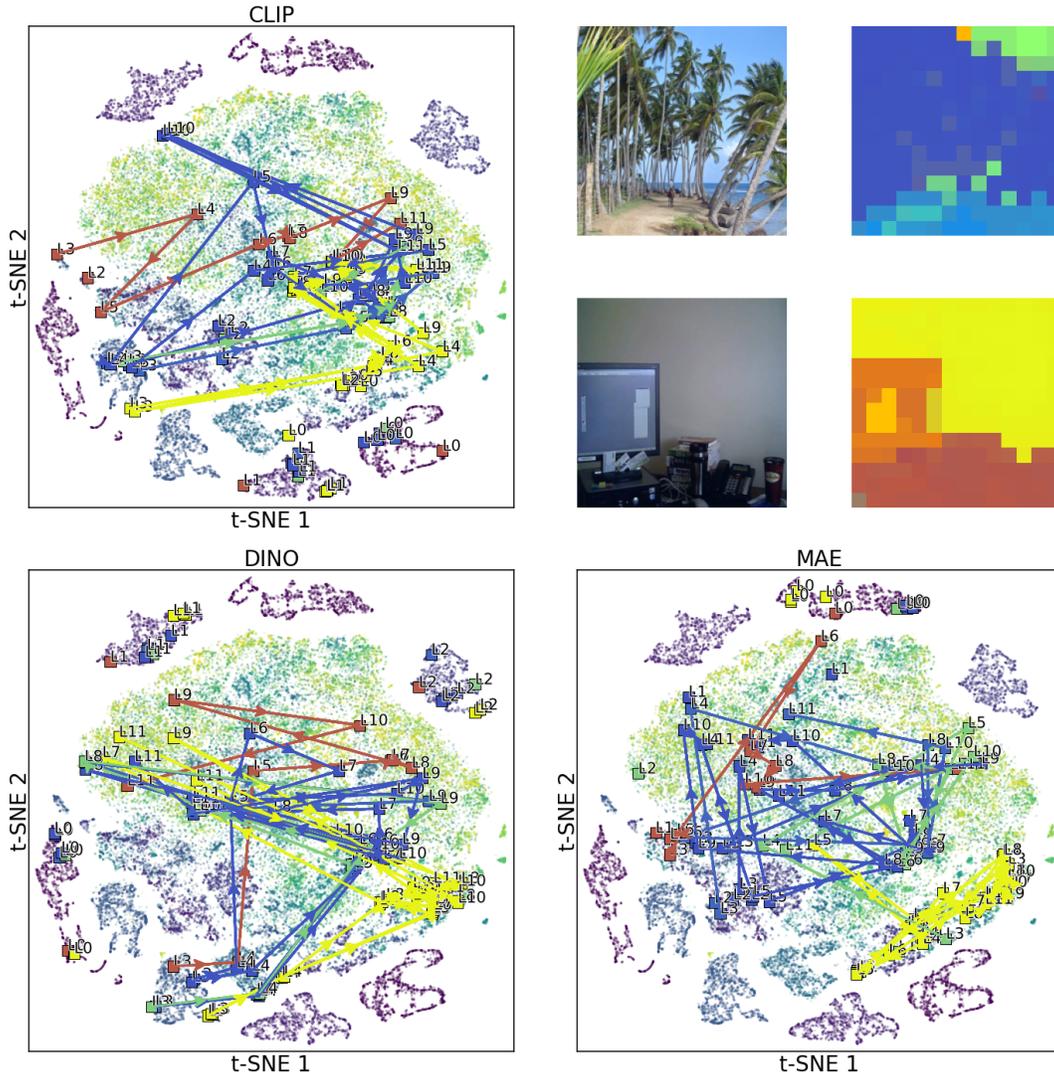


Figure 20: Trajectory of feature progression in from layer to layer, in the 2D spectral-tSNE space. Arrows displayed for 10 randomly sampled example pixels. **Top Right:** Pixels are colored by unsupervised segmentation.

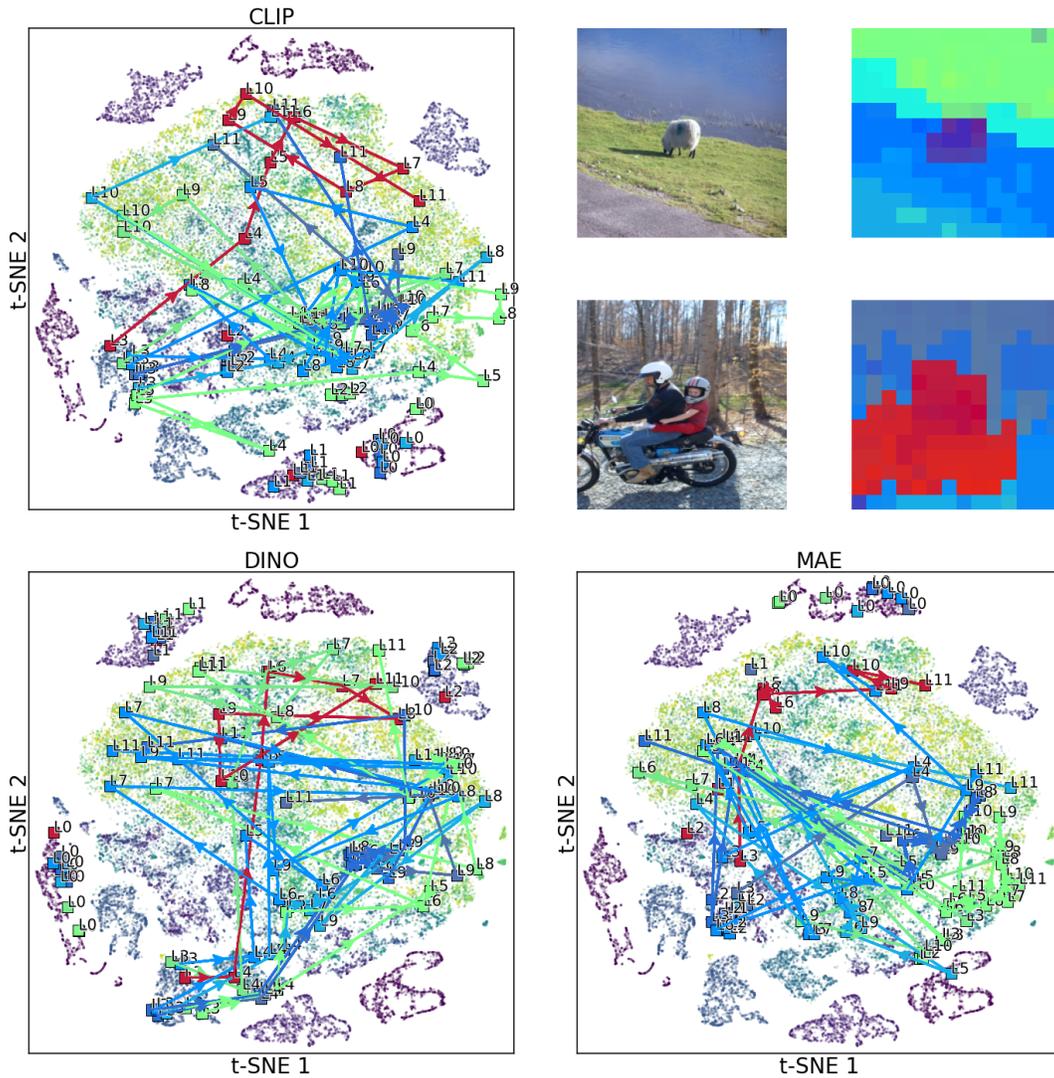


Figure 21: Trajectory of feature progression in from layer to layer, in the 2D spectral-tSNE space. Arrows displayed for 10 randomly sampled example pixels. **Top Right:** Pixels are colored by unsupervised segmentation.

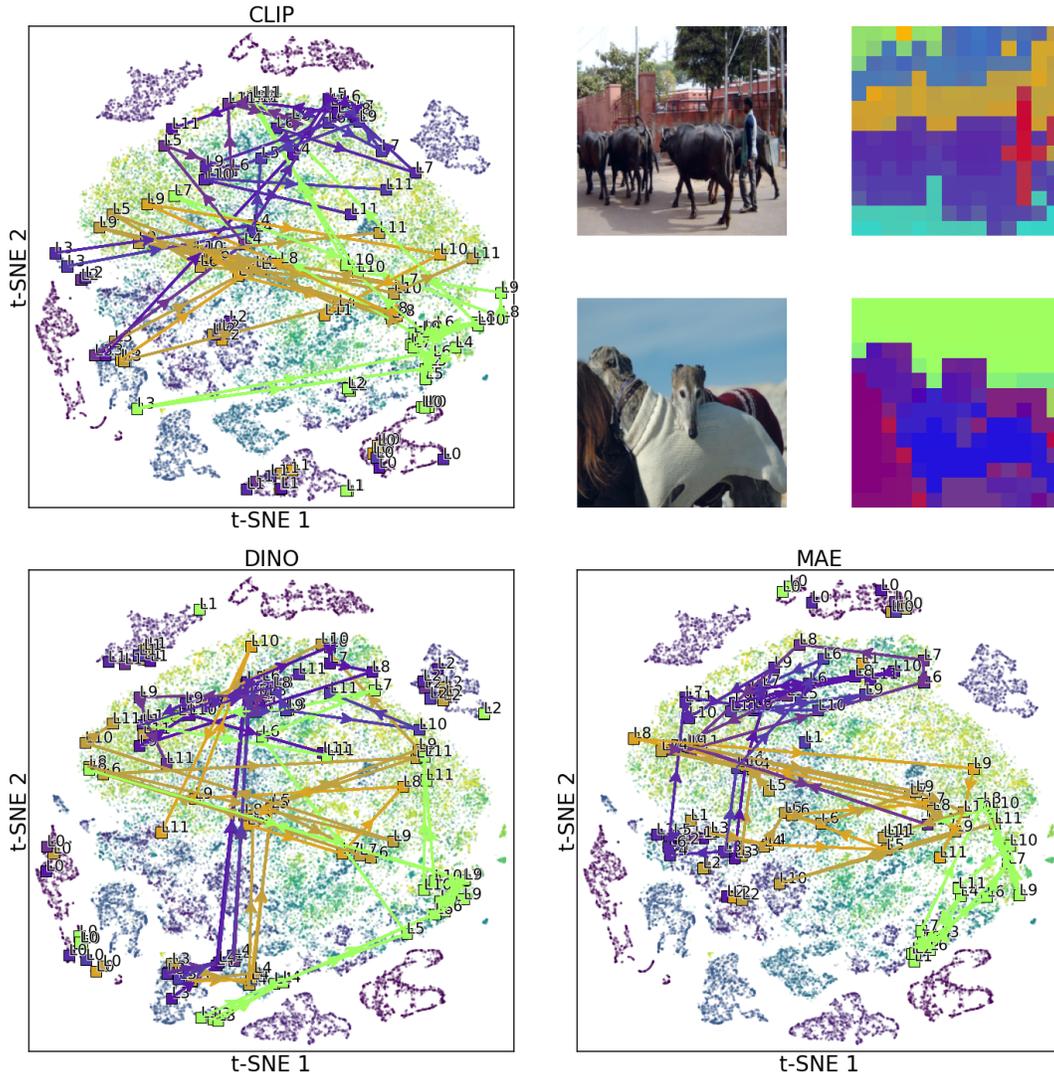


Figure 22: Trajectory of feature progression in from layer to layer, in the 2D spectral-tSNE space. Arrows displayed for 10 randomly sampled example pixels. **Top Right:** Pixels are colored by unsupervised segmentation.

467 **NeurIPS Paper Checklist**

468 **1. Claims**

469 Question: Do the main claims made in the abstract and introduction accurately reflect the
470 paper’s contributions and scope?

471 Answer: [Yes]

472 Justification:

473 Guidelines:

- 474 • The answer NA means that the abstract and introduction do not include the claims
475 made in the paper.
- 476 • The abstract and/or introduction should clearly state the claims made, including the
477 contributions made in the paper and important assumptions and limitations. A No or
478 NA answer to this question will not be perceived well by the reviewers.
- 479 • The claims made should match theoretical and experimental results, and reflect how
480 much the results can be expected to generalize to other settings.
- 481 • It is fine to include aspirational goals as motivation as long as it is clear that these
482 goals are not attained by the paper.

483 **2. Limitations**

484 Question: Does the paper discuss the limitations of the work performed by the authors?

485 Answer: [Yes]

486 Justification:

487 Guidelines:

- 488 • The answer NA means that the paper has no limitation while the answer No means
489 that the paper has limitations, but those are not discussed in the paper.
- 490 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 491 • The paper should point out any strong assumptions and how robust the results are to
492 violations of these assumptions (e.g., independence assumptions, noiseless settings,
493 model well-specification, asymptotic approximations only holding locally). The au-
494 thors should reflect on how these assumptions might be violated in practice and what
495 the implications would be.
- 496 • The authors should reflect on the scope of the claims made, e.g., if the approach was
497 only tested on a few datasets or with a few runs. In general, empirical results often
498 depend on implicit assumptions, which should be articulated.
- 499 • The authors should reflect on the factors that influence the performance of the ap-
500 proach. For example, a facial recognition algorithm may perform poorly when image
501 resolution is low or images are taken in low lighting. Or a speech-to-text system might
502 not be used reliably to provide closed captions for online lectures because it fails to
503 handle technical jargon.
- 504 • The authors should discuss the computational efficiency of the proposed algorithms
505 and how they scale with dataset size.
- 506 • If applicable, the authors should discuss possible limitations of their approach to ad-
507 dress problems of privacy and fairness.
- 508 • While the authors might fear that complete honesty about limitations might be used by
509 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
510 limitations that aren’t acknowledged in the paper. The authors should use their best
511 judgment and recognize that individual actions in favor of transparency play an impor-
512 tant role in developing norms that preserve the integrity of the community. Reviewers
513 will be specifically instructed to not penalize honesty concerning limitations.

514 **3. Theory Assumptions and Proofs**

515 Question: For each theoretical result, does the paper provide the full set of assumptions and
516 a complete (and correct) proof?

517 Answer: [NA]

518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572

Justification:

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: Experimental details are in the appendix

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624

Answer: [No]

Justification: The data is provided as open-source from another study; our code is not yet released, it will be released upon publication.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Experimental details are in the appendix

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We provided standard deviation in Table 1, measured over training with 3 random seed.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer “Yes” if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- 625 • It should be clear whether the error bar is the standard deviation or the standard error
626 of the mean.
- 627 • It is OK to report 1-sigma error bars, but one should state it. The authors should prefer-
628 ably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of
629 Normality of errors is not verified.
- 630 • For asymmetric distributions, the authors should be careful not to show in tables or
631 figures symmetric error bars that would yield results that are out of range (e.g. negative
632 error rates).
- 633 • If error bars are reported in tables or plots, The authors should explain in the text how
634 they were calculated and reference the corresponding figures or tables in the text.

635 8. Experiments Compute Resources

636 Question: For each experiment, does the paper provide sufficient information on the com-
637 puter resources (type of compute workers, memory, time of execution) needed to reproduce
638 the experiments?

639 Answer: [Yes]

640 Justification:

641 Guidelines:

- 642 • The answer NA means that the paper does not include experiments.
- 643 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
644 or cloud provider, including relevant memory and storage.
- 645 • The paper should provide the amount of compute required for each of the individual
646 experimental runs as well as estimate the total compute.
- 647 • The paper should disclose whether the full research project required more compute
648 than the experiments reported in the paper (e.g., preliminary or failed experiments
649 that didn't make it into the paper).

650 9. Code Of Ethics

651 Question: Does the research conducted in the paper conform, in every respect, with the
652 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

653 Answer: [Yes]

654 Justification:

655 Guidelines:

- 656 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 657 • If the authors answer No, they should explain the special circumstances that require a
658 deviation from the Code of Ethics.
- 659 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
660 eration due to laws or regulations in their jurisdiction).

661 10. Broader Impacts

662 Question: Does the paper discuss both potential positive societal impacts and negative
663 societal impacts of the work performed?

664 Answer: [NA]

665 Justification:

666 Guidelines:

- 667 • The answer NA means that there is no societal impact of the work performed.
- 668 • If the authors answer NA or No, they should explain why their work has no societal
669 impact or why the paper does not address societal impact.
- 670 • Examples of negative societal impacts include potential malicious or unintended uses
671 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
672 (e.g., deployment of technologies that could make decisions that unfairly impact spe-
673 cific groups), privacy considerations, and security considerations.

- 674
- 675
- 676
- 677
- 678
- 679
- 680
- 681
- 682
- 683
- 684
- 685
- 686
- 687
- 688
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

689 11. Safeguards

690 Question: Does the paper describe safeguards that have been put in place for responsible
691 release of data or models that have a high risk for misuse (e.g., pretrained language models,
692 image generators, or scraped datasets)?

693 Answer: [NA]

694 Justification:

695 Guidelines:

- 696
- 697
- 698
- 699
- 700
- 701
- 702
- 703
- 704
- 705
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

706 12. Licenses for existing assets

707 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
708 the paper, properly credited and are the license and terms of use explicitly mentioned and
709 properly respected?

710 Answer: [Yes]

711 Justification:

712 Guidelines:

- 713
- 714
- 715
- 716
- 717
- 718
- 719
- 720
- 721
- 722
- 723
- 724
- 725
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

726 • If this information is not available online, the authors are encouraged to reach out to
727 the asset’s creators.

728 **13. New Assets**

729 Question: Are new assets introduced in the paper well documented and is the documenta-
730 tion provided alongside the assets?

731 Answer: [NA]

732 Justification:

733 Guidelines:

- 734 • The answer NA means that the paper does not release new assets.
- 735 • Researchers should communicate the details of the dataset/code/model as part of their
736 submissions via structured templates. This includes details about training, license,
737 limitations, etc.
- 738 • The paper should discuss whether and how consent was obtained from people whose
739 asset is used.
- 740 • At submission time, remember to anonymize your assets (if applicable). You can
741 either create an anonymized URL or include an anonymized zip file.

742 **14. Crowdsourcing and Research with Human Subjects**

743 Question: For crowdsourcing experiments and research with human subjects, does the pa-
744 per include the full text of instructions given to participants and screenshots, if applicable,
745 as well as details about compensation (if any)?

746 Answer: [NA]

747 Justification:

748 Guidelines:

- 749 • The answer NA means that the paper does not involve crowdsourcing nor research
750 with human subjects.
- 751 • Including this information in the supplemental material is fine, but if the main contri-
752 bution of the paper involves human subjects, then as much detail as possible should
753 be included in the main paper.
- 754 • According to the NeurIPS Code of Ethics, workers involved in data collection, cura-
755 tion, or other labor should be paid at least the minimum wage in the country of the
756 data collector.

757 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
758 Subjects**

759 Question: Does the paper describe potential risks incurred by study participants, whether
760 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
761 approvals (or an equivalent approval/review based on the requirements of your country or
762 institution) were obtained?

763 Answer: [NA]

764 Justification:

765 Guidelines:

- 766 • The answer NA means that the paper does not involve crowdsourcing nor research
767 with human subjects.
- 768 • Depending on the country in which research is conducted, IRB approval (or equiva-
769 lent) may be required for any human subjects research. If you obtained IRB approval,
770 you should clearly state this in the paper.
- 771 • We recognize that the procedures for this may vary significantly between institutions
772 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
773 guidelines for their institution.
- 774 • For initial submissions, do not include any information that would break anonymity
775 (if applicable), such as the institution conducting the review.