

# Visual Concept Connectome (VCC): Open World Concept Discovery and their Interlayer Connections in Deep Models

Matthew Kowal<sup>1,3</sup> Richard P. Wildes<sup>1,2</sup> Konstantinos G. Derpanis<sup>1,2,3</sup>

<sup>1</sup>York University, <sup>2</sup>Samsung AI Centre Toronto, <sup>3</sup>Vector Institute

Project Page: [yorkucvil.github.io/VCC](https://yorkucvil.github.io/VCC)

## Abstract

*Understanding what deep network models capture in their learned representations is a fundamental challenge in computer vision. We present a new methodology to understand such vision models, the Visual Concept Connectome (VCC), which discovers human interpretable concepts and their interlayer connections in a fully unsupervised manner. Our approach simultaneously reveals fine-grained concepts at a layer, connection weightings across all layers and is amendable to global analysis of network structure (e.g. branching pattern of hierarchical concept assemblies). Previous work yielded ways to extract interpretable concepts from single layers and examine their impact on classification, but did not afford multilayer concept analysis across an entire network architecture. Quantitative and qualitative empirical results show the effectiveness of VCCs in the domain of image classification. Also, we leverage VCCs for the application of failure mode debugging to reveal where mistakes arise in deep networks.*

## 1. Introduction

This paper focuses on interpreting the intermediate representations of deep networks for computer vision. The goal is understanding how various encoded concepts impact a model’s prediction as well as concepts in other layers. We define concepts as abstractions that generalize from particular instances, including those defined locally (e.g. color and orientation), regionally (e.g. texture and shading) and from higher-level considerations (e.g. object parts, wholes and groupings). Human-interpretable concepts are of particular interest in increasing human understanding of models; however, extracting these concepts encoded in deep networks remains an open challenge in computer vision due to the complexity and opaque nature of these models.

Understanding what concepts are learned by deep models and how they are encoded is important for both science and applications. For science, understanding what information drives a model’s encoding of different concepts may indicate directions to advance performance. For applications,

several negative consequences of deploying opaque vision models have been documented, e.g. [8, 24].

Previous work has focused on interpreting models via feature attribution, which measures the contribution of individual inputs to a model’s output [10, 46, 59]; so, explanations are of single pixels and may be difficult to interpret. Other work generates images that maximize activation of a model’s features [19, 38, 54]. Like feature attribution, these approaches are qualitative and place most of the burden on the user to determine the concepts revealed. Concept-based interpretability, which identifies human-interpretable abstractions in a model’s latent space [20, 23, 29, 32], yield quantitative contributions of a concept to the model’s output and explanations on the class-level (vs. pixel-level). These approaches are easier to understand and validate; however, they have not been used to explore *interlayer* relationships.

As it stands, no approaches can quantify the interlayer affect of a given distributed concept at one layer,  $l$ , to another concept at a different layer,  $l'$  (rather than the model *output*). Even though it is well established that deep networks learn to build concepts hierarchically as information flows through the network [39, 54], understanding the hierarchical representations has been under-researched. Indeed, little is known about the characteristics of this concept hierarchy for today’s models. Questions abound: *How many concepts exist in a network? What are the connections and weights between concepts? Does the model architecture impact the hierarchical structure of concept abstractions?*

In response to these questions, we take inspiration from the biological notion of a *connectome* [47], defined as “a comprehensive structural description of the network of elements and connections forming a brain.” Analogously, we present the *Visual Concept Connectome* (VCC), a comprehensive structural description of a deep pretrained network model in terms of human-interpretable concepts and their relationships that form the internal representation maintained by the model; Fig. 1 shows an example. Notably, VCC generation works in the *open-world* setting, i.e. the concepts and interlayer connections are discovered without the need for any predefined concept dictionary.

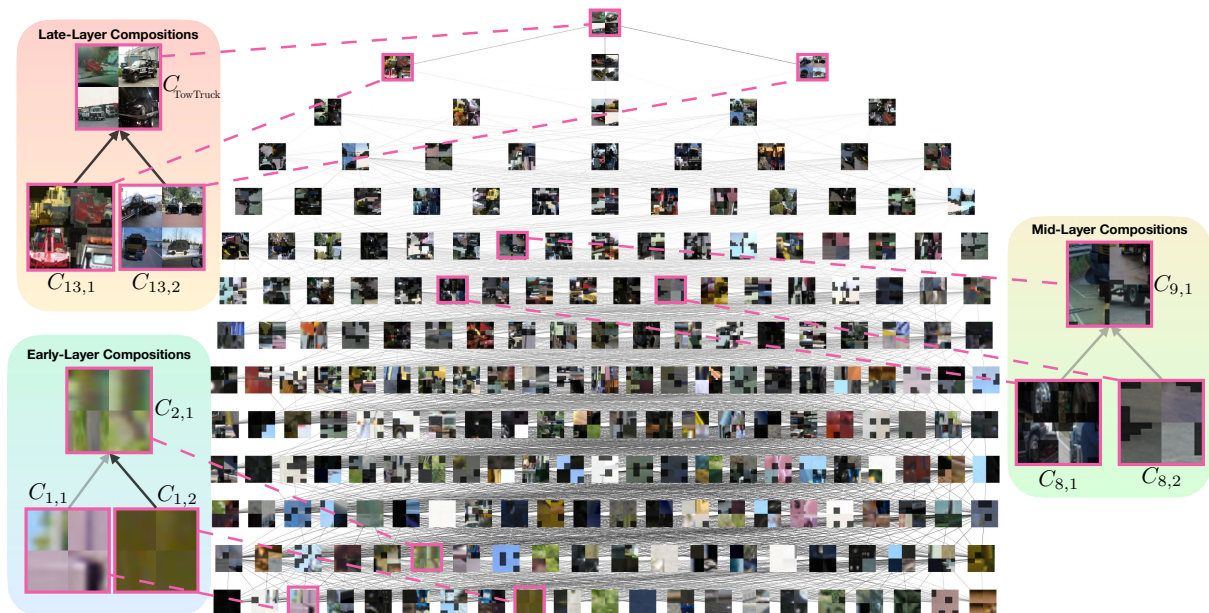


Figure 1. A Visual Concept Connectome (VCC). At each layer the visual concepts learned by a deep model for a given class are revealed as are the learned interlayer concept connections. For each concept, up to four exemplars are shown as unmasked regions in a  $2 \times 2$  image. Interlayer concept connections are shown as lines with darker lines indicating larger contributions. Shown is a VCC for every convolutional layer of a VGG16 model [49] trained on ImageNet [15] targeting recognition of class “Tow Truck”. A closer visualization of VCC subgraphs reveals interesting compositions occurring at different levels of abstraction corresponding at different depths of the model. At early layers (bottom left), we observe oriented patterns ( $C_{1,1}$ ) and brown color ( $C_{1,2}$ ) composing the concept of green and brown orientation ( $C_{2,1}$ ). Middle layers (right) show the concept of ‘wheel on the road’ ( $C_{9,1}$ ) being composed of wheels ( $C_{8,1}$ ) and regions of asphalt ( $C_{8,2}$ ). The final layer concepts (top left) show that both foreground objects, *e.g.* tow trucks ( $C_{13,1}$ ), and background regions, *e.g.* road, trees, humans, or car being towed ( $C_{13,2}$ ), concepts highly influence the final category ( $C_{TowTruck}$ ).

## 2. Related research

While a number of intrinsically interpretable networks have appeared (*e.g.* [6, 12, 31, 45, 56, 57]), we focus on work that, like ours, endeavours to interpret black-box models.

**Concept-based interpretability.** Closed-world concept interpretability considers cases where a labelled dataset defines the concepts of interest for post-hoc analysis; approaches include Network Dissection [3] and Testing with Concept Activation Vectors (TCAV) [29]. However, a desirable property for concept-based interpretability is not to be restricted to a set of predefined classes (*i.e.* closed-world), but to support discovery of new and previously unlabelled concepts (*i.e.* open-world). Approaches to open-world concept interpretability include Automatic Concept Explanations (ACE) [23], SegDiscover [27], ConceptSHAP [53], Invertible Concept Explanations [58] and CRAFT [20]. These approaches are limited to single layer analysis and do not explain interlayer relationships.

Activation maximization [19, 36, 38, 54] visualizes the input that most activates a model component (*e.g.* filter, layer or logit) or how same layer units combine information [40]. These approaches do not capture interlayer relations and place a heavy burden on the user to interpret what concepts appear in their output where there may be no clear resemblance to natural concepts. Class Activation Maps

(CAMs) and extensions visualize an input image’s local region that contributes to a model’s output [9, 46, 59]. Layer-wise Relevance Propagation [5] (LRP) generates pixel contribution heatmaps by assuming conservation of information is propagated through each neuron. These approaches interpret single images (not classes), are purely qualitative, and are sensitive to small input perturbations [22, 30], yet insensitive to model changes [2]. Less related are directions interpreting generative models, *e.g.* [4, 34].

**Interpretability via decomposition.** Most closely related to our work are those that decompose a model’s prediction into interpretable components. CRAFT [20] extracts concepts from the last layer of a convolutional neural network (CNN), but also searches earlier layers for sub-concepts for the maximally activating images of a given concept, but do not provide any way to quantify the contribution of sub-concepts to concepts. Another approach combined CAMs of multiple concepts to explain the final prediction [60]; however, it is limited to labelled data (so cannot discover concepts outside training) and only explains a final output, not intermediate relations. Other work interprets CNNs by distilling them into a graph [55]; however, it applies only to individual filters conceived as capturing object parts (*i.e.* not concepts other than objects and their parts, *e.g.* textures, colors or grouping), does not con-

sider the fact that representations are encoded via superposition [7, 17, 40] (*i.e.* more than one concept can be captured per-channel), does not yield meaningful edge weights and is only for CNNs. Yet other work extends Grad-CAM or LRP to perform intra-layer visualizations [1, 13]. While these approaches capture concepts other than object parts, they still ignore the core problem of superposition (and distributed representations) and do not yield interpretation for an entire class, just single images.

**Contributions.** In the light of previous work, our contributions are fourfold. (i) We introduce and formalize the notion of a Visual Concept Connectome (VCC) for deep network models. The VCC reveals concepts represented at any given layer of a network as well as their interlayer connectivity. (ii) We present a method for extracting a VCC from any pretrained deep network in an unsupervised, open world setting, with a focus on classification. (iii) We validate our approach with quantitative and qualitative experiments wherein we examine various standard models to yield insights into model architectures and training tasks. (iv) We apply VCCs to explaining failure modes in deep models.

### 3. Visual Concept Connectomes (VCCs)

A VCC is a directed acyclic graph,  $\mathbf{G}(\mathbf{Q}, \mathbf{E})$ , created by distilling a pretrained deep network. The graph is topologically sorted based on the layer order found in the original network and has  $n+1$  layers, consisting of  $n$  network layers to be analyzed (*i.e.* any subset of layers that a user selects, including all layers at the extreme) plus the final prediction (*e.g.* the object category for object recognition). The graph’s nodes,  $\mathbf{Q}$ , are vectors (cluster centroids) representing interpretable concepts and its edges,  $\mathbf{E}$ , are scalars representing the contribution of one concept to the existence of another.

To construct a VCC, three inputs are required: (i) a set of  $I$  images representative of a given task (*e.g.* exemplar images of an object category for object recognition),  $\mathcal{I} = \{\mathcal{I}^1, \dots, \mathcal{I}^I\}$ ,  $\mathcal{I}^i \in \mathbb{R}^{h \times w \times c}$ , with  $h, w$  and  $c$  the image height, width and channel dimension (*e.g.* three for RGB), respectively, (ii) an  $N$  layer network,  $F(\cdot) = \{f_1(\cdot), \dots, f_N(\cdot)\}$ , with  $f_j$  denoting feature extraction at layer  $j$ , and (iii) a set of  $n$  selected layers (a subset of  $F$ , possibly improper), which are to be studied.

A VCC is constructed in three main steps. (i) Image segments are extracted in the model’s feature space via divisive clustering to produce semantically meaningful image regions for each selected layer (Sec. 3.1). (ii) Layer-wise concepts, *i.e.* the nodes of the graph, are discovered in an open world fashion (*i.e.* no labelled data is required) via a second round of clustering over the dataset of image regions, independently for each layer (Sec. 3.2). (iii) Edges are calculated that indicate the contribution of concepts from earlier to deeper layers via an approach we introduce, Inter-layer Testing with Concept Activation Vectors (ITCAVs)

(Sec. 3.3). We use object recognition for the explanation of the approach; nevertheless, it could be extended to other tasks in a straightforward way (*e.g.* semantic segmentation). An overview of these steps is provided in Fig. 2.

#### 3.1. Feature space image segments

For each selected layer, we produce a set of image regions that plausibly belong to a concept encoded by the model at that layer. Previous work produces concept segments in RGB space, *i.e.* superpixels or random crops [20, 23]; however, segmentation based on features different from those over which concepts are to be discovered, *i.e.* the model’s deep features, results in the support for the discovery process being divorced from the process that generated the support. Therefore, our segmentation uses the same deep features to be used subsequently for concept discovery. Our divisive clustering approach is presented visually in Fig. 2 (A). We pass the entire set of  $I$  input images,  $\mathcal{I}$ , for a given class through the model,  $F$ , to get features,  $\mathbf{z}_n^i \in \mathbb{R}^{h_n \times w_n \times c_n}$ , at each selected layer,  $n$  according to

$$\mathbf{Z}_n = \{f_n(\mathcal{I}^1), \dots, f_n(\mathcal{I}^I)\} = \{\mathbf{z}_n^1, \dots, \mathbf{z}_n^I\}. \quad (1)$$

To extract image segments at layer  $n$  for concept discovery (Sec. 3.2), we cluster the image activations,  $\mathbf{Z}_n$ , conditioned on the clusters (*i.e.* masks) generated at the next higher layer,  $n+1$ . Let  $\mathbf{p} = (x, y)$  index spatial coordinates and  $\mathbf{B}_n^i(\mathbf{p}; \gamma)$  be a binary mask for cluster  $\gamma$  such that  $\mathbf{B}_n^i(\mathbf{p}; \gamma) = 1 \iff \mathbf{z}_n^i(\mathbf{p}) \in \gamma$ ; otherwise,  $\mathbf{B}_n^i(\mathbf{p}; \gamma) = 0$ .

We calculate a set of such masks by applying a clustering algorithm,  $C_\Gamma^{seg}(\cdot)$ , that returns binary support masks for  $\Gamma$  clusters on its argument. Let  $\Gamma_{n+1}$  be the number of segments at layer  $n+1$  and  $1 \leq g \leq \Gamma_{n+1}$  index a particular segment,  $g$ . For each  $g$  we calculate segments at layer  $n$  as

$$\{\mathbf{B}_n^i(\mathbf{p}; \gamma)\}_{\gamma=1}^{\Gamma} = C_\Gamma^{seg}(\mathbf{z}_n^i(\mathbf{p}) \odot \tilde{\mathbf{B}}_{n+1}^i(\mathbf{p}; g)), \quad (2)$$

where  $\tilde{\mathbf{B}}_{n+1}^i$  is  $\mathbf{B}_{n+1}^i$  upsampled to the resolution of layer  $n$  and  $\odot$  indicates spatial element-wise masking applied to the individual image activations. The element-wise masking ensures clustering is done on the masked area of  $\mathbf{z}_n^i$ . Thus, the support of the concept discovery at layer  $n$  comes from the support of concepts at layer  $n+1$ . Notably, we let  $\Gamma$  vary with  $g$ ; so, different segments,  $g$ , at layer  $n+1$  can yield a different number of segments,  $\Gamma$ , at layer  $n$ .

The top-down conditional clustering, (2), is repeated recursively for all image segments in all selected layers. To initiate the recursion, (2), at the top layer,  $n_{top}$ , all the features are used by setting the number of clusters to be one, and having  $\tilde{\mathbf{B}}_{n_{top}+1}^i(\mathbf{p}; 1) = 1$  for all images.

At each layer,  $n$ , we construct a set of RGB segmentations,  $\mathbf{M}$ , for use in concept discovery (Sec. 3.2) by applying upsampled masks,  $\mathbf{B}_n^i(\mathbf{p}; \gamma)$ , to a given image,  $\mathcal{I}^i$ , via

$$\mathbf{M}_n^i(\mathbf{p}; \gamma) = \mathcal{I}^i(\mathbf{p}) \odot \mathbf{B}_n^i(\mathbf{p}; \gamma), \quad (3)$$

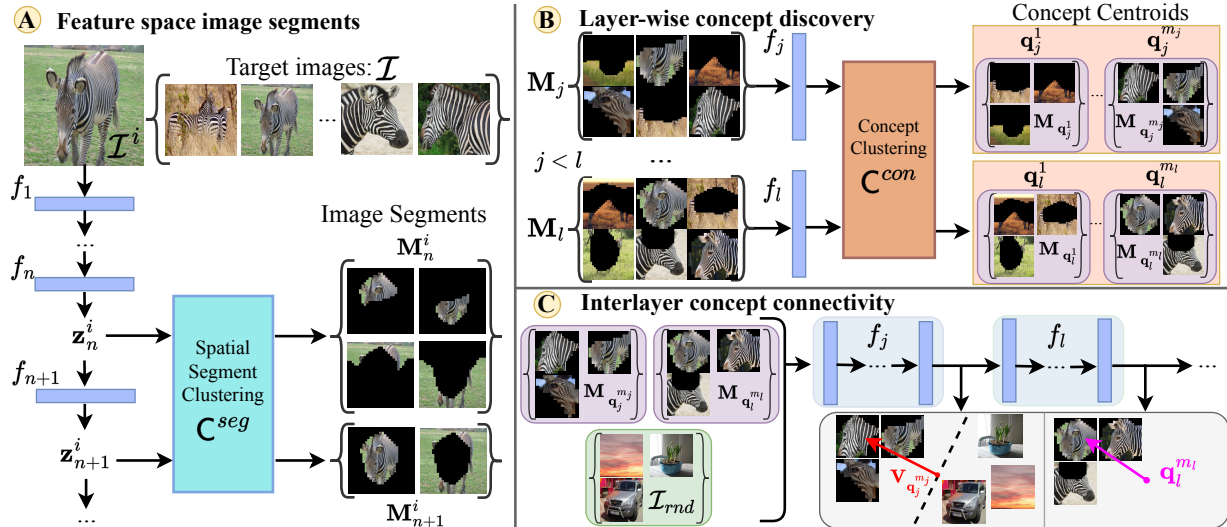


Figure 2. The three steps in building a Visual Concept Connectome (VCC). (A) For a given image,  $\mathcal{I}^i \in \mathcal{I}$ , model,  $F$ , and layer,  $n$ , we produce a set of image segments,  $\mathbf{M}_n^i \in \mathbf{M}_n$ , based on a recursive spatial clustering,  $\mathbf{C}^{seg}$  (2), of the features  $\mathbf{z}_n^i$  conditioned on the clusters from the layer above,  $n+1$ . We then use (3) to generate a set of masked RGB image segments for each layer,  $\mathbf{M}_j$ . (B) For a given layer,  $j$ , we pass the image segments from all images,  $\mathbf{M}_j$ , through  $f_j$  and cluster,  $\mathbf{C}^{con}$  (4), these features across the dataset to produce  $m_j$  concept centroids,  $\{\mathbf{q}_j^1, \dots, \mathbf{q}_j^{m_j}\}$ . (C) To measure the contribution of an earlier layer concept,  $\mathbf{q}_j^{m_j}$ , to a later layer concept,  $\mathbf{q}_l^{m_l}$ , we employ our Interlayer Testing with CAV (ITCAV) approach (Sec. 3.3), which uses the Concept Activation Vector (CAV) [29] of the earlier concept,  $\mathbf{V}_{\mathbf{q}_j^{m_j}}$  (that points away from random examples,  $\mathcal{I}_{rnd}$ , but toward concept exemplars,  $\mathbf{M}_{\mathbf{q}_j^{m_j}}$ ), and the deeper layer concept,  $\mathbf{q}_l^{m_l}$ .

where  $\mathbf{M}_n^i(\mathbf{p}; \gamma) \in \{0, \dots, 255\}^{h_n \times w_n}$ , with  $\{0, \dots, 255\}$  specifying RGB value. We follow by defining  $\mathbf{M}_n$  as the set of all RGB image segments at layer  $n$  (with each segment given in terms of (3)), and letting  $\mathbf{M} = \{\mathbf{M}_1, \dots, \mathbf{M}_n\}$ .

We instantiate the clustering,  $\mathbf{C}^{seg}$ , via maskSLIC [28], an extension of k-means, with an  $l_2$  norm, that clusters features while respecting a mask and automate selection of the number of clusters via the silhouette method [43].

### 3.2. Layer-wise concept discovery

Given the dataset of masked image segments for each layer,  $\mathbf{M} = \{\mathbf{M}_1, \dots, \mathbf{M}_n\}$ , we follow [23] and bilinearly interpolate the segments to the input image size and then pass the segments through the model to get pre-segmented activations  $\mathbf{Z}_{\mathbf{M}_j} = f_j(\mathbf{M}_j)$  from the layer where the segments were found in the divisive clustering step (Sec. 3.1). We pass the image segments (as opposed to using the image activations) through the model to remove any contextual information before we discover concepts (the segments are set to zero everywhere outside of the segment boundary). Otherwise,  $\mathbf{M}_j$  could mix information from outside the segment.

The layer-wise concept discovery step is presented visually in Fig. 2 (B). With the pre-segmented activations,  $\mathbf{Z}_{\mathbf{M}_j}$ , calculated for all layers,  $j$ , we cluster over the dataset of segment features for each individual layer to produce cluster centroids (*i.e.* concepts) according to

$$\mathbf{Q}_j = \{\mathbf{q}_j^1, \dots, \mathbf{q}_j^{m_j}\} = \mathbf{C}_{m_j}^{con}(\text{GAP}(\mathbf{Z}_{\mathbf{M}_j})), \quad (4)$$

where  $m_j$  is the number of concepts discovered at layer  $j$ ,

GAP is spatial global average pooling and  $\mathbf{C}_{m_j}^{con}$  is a clustering algorithm that returns  $m_j$  cluster centroids. Following previous work [23], we instantiate  $\mathbf{C}_{m_j}^{con}$  in terms of standard k-means, using the  $l_2$  norm, and then use a large number of clusters with subsequent pruning to remove noisy clusters; see the appendix for details. The output of this stage results in discovered concepts (*i.e.* cluster centroids and associated segments) for all layers as  $\mathbf{Q} = \{\mathbf{Q}_1, \dots, \mathbf{Q}_n\}$  for the centroids and  $\mathbf{M}_{\mathbf{q}_j^{m_j}}$  their associated RGB segments.

### 3.3. Interlayer concept connectivity

Our approach to quantifying the contribution of discovered concepts between two layers, Interlayer Testing with Concept Activation Vectors (ITCAV), generalizes TCAV [29]. Specifically, while TCAV calculates the contribution of a single mid-layer concept to the class logit using a dataset of labelled concept images, we generalize to operate between any two layers and without a dataset of labelled concept images. Without loss of generality, we consider a single pair of concepts at two layers, which need not be adjacent.

The ITCAV method is presented visually in Fig. 2 (C). We seek to measure the degree to which a concept at an earlier layer contributes to a deeper layer. To do so, we construct two vectors: (i) a vector that represents the concept at the earlier layer,  $j$ , and (ii) a vector that represents a positive contribution (*i.e.* gradient) of a concept at a later layer,  $l$ . The closer the alignment of these two vectors, the larger the contribution of the concept at layer  $j$  to the concept at layer  $l$ . Let  $\mathbf{q}_j^{m_j}$  and  $\mathbf{M}_{\mathbf{q}_j^{m_j}}$  be the cluster centroid and as-

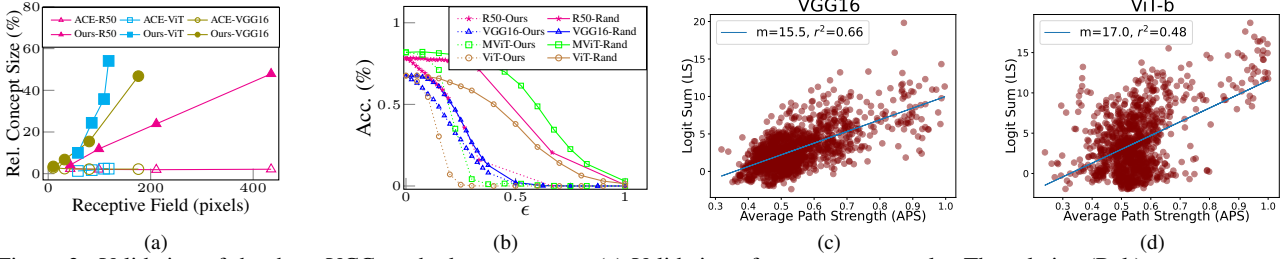


Figure 3. Validation of the three VCC method components. (a) Validation of segment proposals. The relative (Rel.) concept segment size compared to entire image for a given layer is plotted against the receptive field (RF) width/height of the same layer. (b) Validation of discovered concepts. For 50 randomly selected ImageNet classes, we discover concepts in four layers of the model. During inference, one randomly selected concept at each layer is suppressed by a factor of  $\epsilon$ . (c, d) Validation of interlayer concept weights. The unnormalized logit sum (LS) scores, (8), for the target class are plotted against the average path strength (APS) scores, (7). A positive correlation implies that the ITCAV edge weights connecting a concept to the class are predictive of the model output having a higher probability for that class.

sociated RGB image segments in layer  $j$ , and let  $\mathbf{q}_l^{m_l}$  and  $\mathbf{M}_{\mathbf{q}_l^{m_l}}$  be the cluster centroid and associated RGB image segments in a deeper layer  $l$ , *i.e.*  $l > j$ .

We proceed by constructing the vector that represents the concept in the feature space of the earlier layer,  $j$ . To do so, we employ the Concept Activation Vector (CAV) [29] for the earlier concept,  $\mathbf{q}_j^{m_j}$ , by training a linear classifier,  $h(\cdot)$ , on the features of layer  $j$  between positive concept inputs,  $f_j(\mathbf{M}_{\mathbf{q}_j^{m_j}})$ , and random images,  $f_j(\mathcal{I}_{rnd})$ . The orthogonal vector to the hyperplane defined by  $h(\cdot)$ , denoted as  $\mathbf{V}_{\mathbf{q}_j^{m_j}}$ , is the CAV for concept  $\mathbf{q}_j^{m_j}$  and points in the direction of concept  $\mathbf{q}_j^{m_j}$  in the feature space of layer  $j$ .

Next, we construct the vector in the feature space of layer  $j$  that points in the direction of positive contribution of concept  $\mathbf{q}_l^{m_l}$ . To do so, we calculate  $\nabla_{f_l}(f_j(\cdot))$ , *i.e.* the gradient of the deeper concept,  $\mathbf{q}_l^{m_l}$ , at layer  $l$ , with respect to the earlier layer,  $j$ . Our approach to measuring the sensitivity between the two concepts is then

$$S_{\mathbf{q}_j^{m_j}, \mathbf{q}_l^{m_l}}(x) = \nabla \left( \|f_l(f_j(x)) - \mathbf{q}_l^{m_l}\|_2 \right) \cdot \mathbf{V}_{\mathbf{q}_j^{m_j}}, \quad (5)$$

where  $x \in \mathbf{M}_{\mathbf{q}_l^{m_l}}$ . The sensitivity score should be positive because when the feature at layer  $j$  is perturbed in the direction of concept  $\mathbf{q}_j^{m_j}$ , *i.e.* in the direction  $\mathbf{V}_{\mathbf{q}_j^{m_j}}$ , then the feature vector in deeper layer,  $l$ , will be pushed closer to concept  $\mathbf{q}_l^{m_l}$ . Following TCAV [29], the final ITCAV edge weight,  $e_{\mathbf{q}_j^{m_j}, \mathbf{q}_l^{m_l}} \in \mathbf{E}$ , is the ratio of positive sensitivities,

$$e_{\mathbf{q}_j^{m_j}, \mathbf{q}_l^{m_l}} = \left| x \in \mathbf{M}_{\mathbf{q}_l^{m_l}} : S_{\mathbf{q}_j^{m_j}, \mathbf{q}_l^{m_l}}(x) > 0 \right| / \left| \mathbf{M}_{\mathbf{q}_l^{m_l}} \right|. \quad (6)$$

Key to our sensitivity score, (5), is the  $l_2$  distance between the output,  $f_l(f_j(x))$ , and cluster centroid,  $\mathbf{q}_l^{m_l}$ , before taking the gradient, *cf.* CAV [29]. This choice serves two goals: (i) The subset of dimensions of individual image features,  $f_j(x)$ , that are not well aligned with the cluster centroid,  $\mathbf{q}_l^{m_l}$ , will be penalized. This penalty will suppress the impact of noisy dimensions from individual samples on ITCAV. We use the  $l_2$  norm because the employed clustering approach used during concept discovery, k-means, uses

$l_2$ ; see Sec. 3.2. (ii) Collapsing the tensor to a scalar (*i.e.* a concept similarity score) makes the gradient calculation computationally feasible compared to the full Jacobian.

## 4. Experiments

### 4.1. Implementation details

For evaluation we use ResNet50 [25], VGG16 [48], MobileNetV3 [26], ViT-b [51] and MViT [18] to sample a broad range of popular architectures. The exact layers selected are in the appendix, but always include the last feature layer before the fully connected layer to probe concepts closest to the model output. Following previous work [23, 29], we perform a two-sided t-test on ITCAV scores for multiple runs of the same concept vs. different random sets of images to remove statistically insignificant scores and randomly sample from the Broden dataset [3] to get the images,  $\mathcal{I}_{rnd}$  (Sec. 3.3), used in statistical testing and CAV training. On average, four and all-layer VCCs take 15 minutes and 36 hours, resp., to generate on an NVIDIA Quadro RTX 6000 GPU. The appendix has more details and VCCs for other layers, models and datasets [52].

### 4.2. VCC component validation

**Segment proposal validation.** Our approach segments concepts based on the feature space of a given layer to allow us to capture concepts across multiple layers (Sec. 3.1). So, the spatial support of valid segments at a layer should follow the receptive field (RF) of filters at that layer, as the segments are constrained by the filters from which they are derived. For comparison, we consider ACE [23], which does segmentation with a layer independent feature (color), limited to the input image to discover concepts at a single layer for a single model. Being directly related to the features at a layer, our approach should produce concepts that better respect the RF at that layer compared to the baseline, whose features are layer independent (*e.g.* at CNN early layers the RF is small; so, intuitively the patch proposals should be smaller as well and conversely for the later layers).

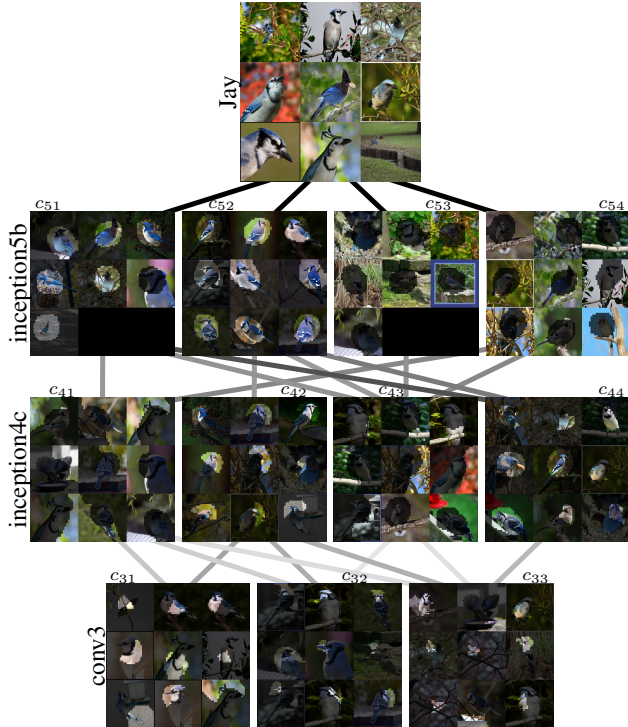


Figure 4. A VCC for three selected layers of a GoogLeNet model [49] targeting recognition of class ‘Jay’. Darker lines denote stronger connection weights.

Figure 3 (a) shows the average concept size (in terms of the ratio of segment pixels to the entire image) vs. the RF at the corresponding layer for three different models: (i) ResNet50 [25], (ii) VGG16 [48] and (iii) ViT-b [16] (note we use mean attention distance instead of receptive field for ViT-b). While the baseline [23] produces segments with sizes invariant to the layer analyzed, our approach produces segments with image sizes that scale with the RF. Additional results with other models (including more transformers) are given in the appendix, which further support this observation. Overall, it is expected by design and we find that our segment proposals follow RFs at each layer.

**Concept fidelity.** Concept fidelity measures the meaningfulness of discovered concepts with respect to the target model. While other work focuses on single layer concept fidelity [20, 23, 29], we measure the effect on a model’s output when suppressing the encoding of concepts as the information propagates through the model. The expectation is that suppressing the discovered concepts should result in a quicker decrease in performance compared to a non-concept direction (*e.g.* a randomly selected direction) as the amount of suppression gets larger. For a given model and category, we compute the concepts at various layers using our method (Sec. 3.2). We perform concept suppression for a given image,  $\mathcal{I}^i$ , at a given layer,  $j$ , for concept  $\mathbf{q}_j$ , according to  $\mathbf{z}_j^i = \mathbf{z}_j^i - \epsilon \cdot \mathbf{q}_j / \|\mathbf{q}_j\|_2$ , where  $\epsilon$  controls the degree of perturbation. If the concepts discovered are meaningful, then the accuracy for the model’s target class should decrease faster

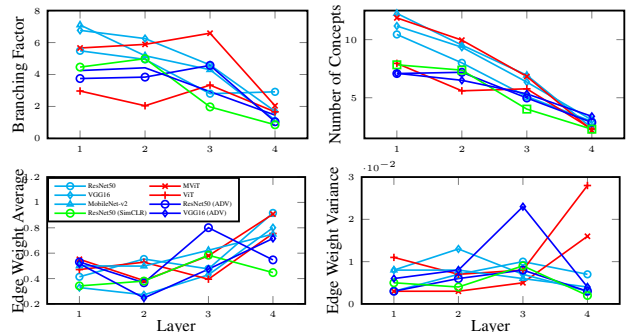


Figure 5. Graph metrics on four layer VCCs comparing CNN vs. transformer architectures and training objectives.

compared to random perturbations,  $\mathbf{q}_{rnd}$ , as  $\epsilon$  increases.

Figure 3 (b) has concept suppression results averaged over 50 randomly selected ImageNet [15] classes (documented in the appendix), where a concept is randomly chosen to be suppressed at each layer. Note that the  $\epsilon$  values required to reduce the model accuracy to zero differs by model; so, we scale the  $\epsilon$  values to  $[0, 1]$  for visualization purposes. For all models, it is seen that perturbing in the opposite direction of a recovered concept more severely impacts the accuracy than a random perturbation.

**ITCAV validation.** If the interlayer concept connection weights are meaningful, then the accumulated weights from earlier to later layer concepts should be correlated with their probability of predicting the target class. We use this intuition to validate our ITC AV approach to recovering interlayer concept weights as follows. We define the average path strength (APS), a scalar value between zero and one representing the average strength of the connection between a concept and the class logit. Given concept  $\mathbf{q}_j$  at layer  $j$ , let there be  $L$  layers in the VCC between the concept and the logit layer; thus, each path from the concept to the logit consists of  $L$  edge weights. We consider all possible paths in the VCC from this concept to the class logit. Let  $e_{\mathbf{q}_j}(l, p)$  be the edge weight (*i.e.* ITC AV score) of the  $p^{\text{th}}$  path (from concept  $\mathbf{q}_j$  to the class concept) for VCC layer  $l$ . To calculate the APS score for a given concept, we average over all paths,  $P$ , and edge weights in each path to get

$$\text{APS}(\mathbf{q}_j) = \frac{1}{P} \sum_{p=1}^P \left[ \frac{1}{L} \sum_{l=1}^L e_{\mathbf{q}_j}(l, p) \right]. \quad (7)$$

Next, we calculate the logit sum (LS) score for a given concept,  $\mathbf{q}_j$ , by passing all of the corresponding masked segments,  $\mathbf{M}_{\mathbf{q}_j}$ , through the model,  $F(\cdot)$ . We then sum up the logit scores for the target class,  $c$ , to yield

$$\text{LS}(\mathbf{q}_j) = \sum_{i=1}^{|\mathbf{M}_{\mathbf{q}_j}|} F(\mathbf{M}_{\mathbf{q}_j}^i)|_c, \quad (8)$$

where  $\mathbf{M}_{\mathbf{q}_j}^i$  denotes the  $i^{\text{th}}$  segment mask, and  $F(\cdot)|_c$  denotes the  $c^{\text{th}}$  logit score. If the ITC AV edge weights are

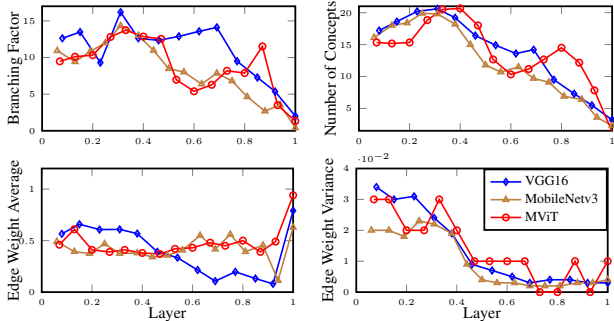


Figure 6. Graph metrics of all layer VCCs for three diverse architectures. Layer number normalized to allow for comparison of models with different numbers of layers.

meaningful, then there should be positive correlation between the APS and LS scores, *i.e.*  $\text{APS}(\mathbf{q}) \propto \text{LS}(\mathbf{q})$ .

Figure 3 (c, d) shows the LS scores plotted vs. the APS scores for VGG16 [25] and ViT-b [51]. A positive correlation is found between the APS and LS scores for both architectures. These results suggest that the combination of ITCAV scores is predictive of whether a concept is representative of the target class. The appendix has more LS vs. APS plots for more models that support these findings.

### 4.3. Understanding models with VCCs

To demonstrate VCC’s unique ability to interpret models at different resolutions, we present model analyses in three parts: (i) qualitative analyses of VCCs generated for a subset of layers as well as at all layers, (ii) quantitative analysis of VCCs generated for a subset of four layers, (iii) a broadened quantitative analysis to encompass *all* model layers.

**Visualizing concept hierarchies.** Figure 4 shows a three layer VCC for GoogLeNet [49] targeting the class “Jay”. Notice how the inception5b ‘bird’ concepts ( $c_{51}, c_{52}$ ) form as selective weighting of inception4c concepts background ( $c_{41}$ ), bird part ( $c_{42}, c_{44}$ ) and tree branch ( $c_{43}$ ), while the inception5b background concepts ( $c_{53}, c_{54}$ ) form differently from weighting of solely inception4c background part concepts (*e.g.* tree branch ( $c_{43}$ ) and green leaves ( $c_{41}$ )). Notably, the network separates subspecies of Jay in the final layer (*e.g.* Blue Jay ( $c_{52}$ ) and other types ( $c_{51}$ )). The concepts found in inception4c are composed from varying combinations of colors and parts found in conv3 (*e.g.* various bird parts ( $c_{31}, c_{33}$ ) contribute to the bird concepts at inception4c). In the end, both scene and object contribute with strong weights to the final category.

An all layer VGG-16 VCC visualization was presented in Fig. 1. As discussed in the caption, hierarchical concept assemblies again are revealed, with both target object as well as its background contributing to the final classification. While both the few layer and all layer visualizations reveal the concept representations of the models under analysis, they afford different levels of granularity: Few layer visualization provides a concise summary with a focus on

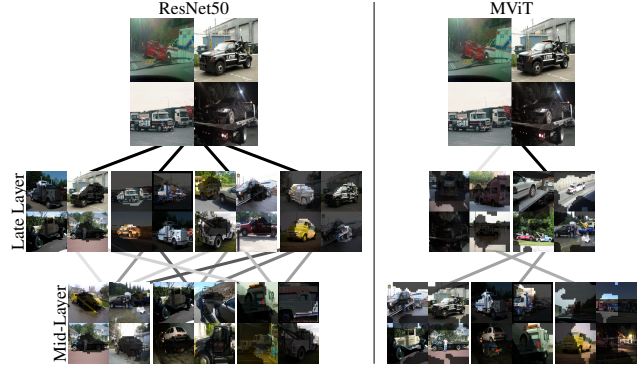


Figure 7. VCCs for a CNN (ResNet50) and transformer (MVIT) for the class ‘Tow Truck’. Shown are only the last two layers with the class output.

specific layers, while all layer provides very detailed study.

**Subset VCC analysis.** We begin our quantitative analysis on four-layer VCCs generated for the CNNs and transformers listed in Sec. 4.1. Layers are selected approximately uniformly (see appendix for exact layers) across a model to capture concepts at different abstractions. We use standard metrics for analyzing tree-like data structures in a per-layer fashion: branching factor (extent to which the representation is distributed), number of nodes (how many concepts the model uses to encode a particular class), edge weights (strength of connection between concepts) and edge weight variance (variability of connection strength). We calculate averages over 50 VCCs for 50 randomly selected ImageNet classes. Results are shown in Fig. 5.

Patterns are apparent in concept numbers and edge weights: more, but weaker, concepts in early layers vs. fewer, but stronger, concepts in later layers. These patterns reflect the shared low-level features (*e.g.* colors, textures) across classes and more specific features near the end, which yield larger ITCAV values. Also, CNNs show a decreasing branching factor, while transformers maintain a consistent number until the final layer, where all models converge to about two concepts, typically of an ImageNet class’s foreground-background structure. Transformers have higher final layer edge weight variance compared to CNNs, indicating their ability to better differentiate earlier concepts’ importance in forming later concepts, potentially explaining their superior classification performance (*i.e.* all information is not equally valuable). We also compare models (ResNet50 and VGG16) when trained with self-supervision [11] (SimCLR) or for adversarial robustness, *i.e.* on Stylized ImageNet [21] (ADV). We observe that robust and self-supervised models have fewer low-level concepts and compositionality than the originals, likely as their training yields less reliance on texture (stylization perturbs texture) and color (SimCLR training jitters color).

To examine these patterns further, VCC visualizations for a CNN and transformer are shown in Fig. 7. Here, we limit to two (diverse) models and two later layers with class

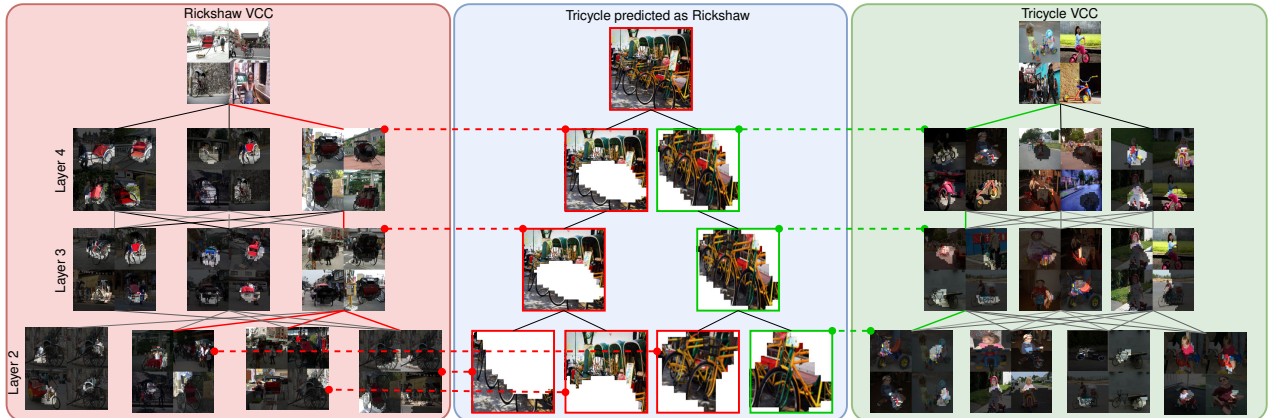


Figure 8. Debugging model failure modes with VCCs. We show an image of a tricycle incorrectly predicted by a ResNet50 as a rickshaw (middle) as well as the top-down segmentation of the image (Sec. 3.1). We also show the incorrect (left) and correct (right) VCCs. Following the hierarchy of concepts reveals that the model incorrectly focused heavily on the collapsible hood, starting at Layer 2.

output for space. The connection diversity in the last layer is indeed observed to be larger for the transformer vs. the CNN. Notably, over half of the concepts in the CNN capture background centric concepts, while the transformer has only a single background centric concept.

**All layer VCCs.** We present quantitative analyses on all layer VCCs for three diverse models: (i) a standard CNN, VGG-16, (ii) an efficient model, MobileNetv3 and (iii) a transformer, MViT. Averages are taken over 10 VCCs for 10 random ImageNet classes. Results are shown in Fig. 6. Common trends appear in all models. Concept composition is non-linear across layers, with branching factor ranging from 5-15 and converging to approximately two near the last layers. The peak number of concepts, around 20, is consistently at 30-40% of network depth and, as in the four-layer analysis, also converges to two in the final layer.

Edge weights and variances are in accord with the main findings of the four-layer analysis, but also reveal other insights into the compositional nature of the models. At fine grained, all layer analysis, each model more readily displays unique edge weight characteristics: VGG16’s average weights decrease in later layers, MobileNetv3’s drop greatly before the final layer and MViT maintains consistent values. Still, overall these results indicate that penultimate concepts differ between CNNs vs. transformers, as supported by our four-layer VCC analysis; see Fig. 5. Higher variances in initial layers suggest a diverse combination of concepts, whereas deeper layers indicate a more uniform composition. Transformers, however, show a variance increase in the final layer, indicating greater compositionality.

In summary, the four and all-layer VCC analyses consistently highlight key aspects of deep networks: (i) Early and mid-layers use more concepts compared to deeper layers. (ii) Concept interactions, *e.g.* ITCAV values and variances, differ greatly between adjacent layers, as opposed to across multiple layers. (iii) CNNs and transformers compose concepts differently, especially at the final layers.

## 5. Application: Diagnosing failure predictions

To show the VCC’s practical utility, we consider application to model failure analysis. Contrasting failure analysis methods that do not explain how errors are made compositionally [46] or only offer explanations of single neuron heatmaps [13], VCCs provide insights on compositional concepts across layers and distributed representations.

Figure 8 shows a ‘Tricycle’ incorrectly classified as a ‘Rickshaw’ by a ResNet50 model, and the corresponding incorrect VCC (‘Rickshaw’, left) and correct VCC (‘Tricycle’, right). As the image is decomposed using our top-down segmentation (Sec. 3.1), it is revealed that the majority of pooled segments are closer, in terms of  $l_2$  distance, to concepts in the Rickshaw VCC (red outlines) than the tricycle VCC (green outlines). While the model correctly encoded the wheel and handlebar portions of the images as tricycle concepts, the background and collapsible hood concepts are composed from layers two through four as rickshaw concepts, which may cause the error. We also note the lack of other tricycle-specific concepts (*e.g.* children).

## 6. Discussion and conclusion

The Visual Concept Connectome (VCC) is a method for discovering human interpretable concepts and their interlayer connections in a fully unsupervised way. VCCs allow for the interpretation of deep models at multiple resolutions, revealing interesting properties via application to a range of networks. When comparing architecture types, CNNs were found to rely heavily on final layer concepts for prediction, while transformers were found to have more diverse connection strengths in the final layer and a more complex assembly of later concepts from earlier ones. We also applied VCCs to model debugging and showed how the hierarchical concept representation can indicate how mistakes arise. Given VCCs generality, our methodology should be applicable to additional networks, tasks and applications.



# Visual Concept Connectome (VCC): Open World Concept Discovery and their Interlayer Connections in Deep Models

## Appendix

### 7. Introduction

This document provides additional material that is supplemental to our main submission. Section 8 outlines the algorithms used in our technical approach. Section 9 describes additional implementation details for our approach, including further Visual Concept Connectome (VCC) generation settings, model details, clustering details and target classes chosen for evaluation. Section 10 provides additional empirical results in terms of validation of the segment proposal method, validation of the concept discovery method, validation of the interlayer testing with concept activation vector (ITCAV) method, VCC visualizations comparing models, classes, and layers, as well as VCCs with a larger number of layers, including all layers. Section 12 discusses the limitations of VCCs and our associated methodology to generate them. Section 13 discusses the societal implications, both positive and negative, of our method. Finally, Section 14 details the used assets and accompanying licenses.

### 8. Algorithms

In this section, we present pseudocode for the three main algorithmic components of our method: (i) Top-down feature segmentation (Sec. 3.1 in the main paper), (ii) Layer-wise concept discovery (Sec. 3.2 in the main paper) and (iii) Interlayer testing with concept activation vectors (ITCAV) (Sec. 3.3 in the main paper). The top-down feature segmentation method is shown in Algorithm 1. The layer-wise concept discovery method is shown in Algorithm 2. Finally, The ITCAV method is shown in Algorithm 3. All references to equations in the algorithms refer to equations in the main paper.

### 9. Implementation details

**VCC settings.** 50 target images are used to generate each VCC. The statistical testing and training of the CAVs [29] use 20 unique sets of random images from the Broden dataset [3]. When computing the concept connection edge weight between the final selected layer and the class logit, the standard TCAV [29] score is used.

**Model settings.** For the CLIP-ResNet50 [42] experiments (Sec. 10.3), we follow the original paper [42] and compute the logit for each ImageNet [15] class using a single query sentence ‘a photo of a {class}’; where ‘{class}’ is the target ImageNet class. The layer names used (according to the PyTorch [41] module nomenclature) for each model when generating all four-layer VCCs are as follows:

---

#### Algorithm 1 Top-Down Feature Segmentation

---

**Input:** Model  $F$ , Set of images  $\mathcal{I}$ ,  $n$  selected layers of  $F$  to study, Spatial clustering algorithm  $C^{seg}$   
*/\* $C^{seg}$  is instantiated in terms of maskSLIC [28]\*/*  
**Output:** Set of RGB image masks  $\mathbf{M}$   
*/\*Set lists for collection of masks and activations\*/*

- 1:  $\mathbf{M} \leftarrow []$ ,
- 2: **for**  $j \in n$  **do**
- 3:      $\mathbf{M}_j, \mathbf{Z}_j \leftarrow [], []$
- 4: **end for**  
*/\*Collect activations at each layer, Eq. (1)\*/*
- 5: **for**  $\mathcal{I}^i \in \mathcal{I}$  **do**
- 6:     **for**  $j \in n$  **do**
- 7:          $\mathbf{z}_j^i \leftarrow f_j(\mathcal{I}^i)$
- 8:          $\mathbf{Z}_j.append(\mathbf{z}_j^i)$
- 9:     **end for**
- 10: **end for**  
*/\*Iterate through all images\*/*
- 11: **for**  $i \in |\mathcal{I}|$  **do**  
*/\*Iterate through layers top-down\*/*
- 12:     **for**  $j \in \{n, \dots, 1\}$  **do**
- 13:          $\mathbf{B}_j^i \leftarrow []$   
*/\*All features considered at top layer\*/*
- 14:         **if**  $j == n$  **then**
- 15:              $\tilde{\mathbf{B}}_{j+1}^i(\mathbf{p}; 1) \leftarrow \{1\}^{h_j \times w_j}$
- 16:              $\Gamma_j^i \leftarrow \text{silhouette}(\mathbf{z}_j^i(\mathbf{p}) \odot \tilde{\mathbf{B}}_{j+1}^i(\mathbf{p}; 1))$
- 17:              $\{\mathbf{B}_j^i(\mathbf{p}; \gamma)\}_{\gamma}^{\Gamma_j^i} \leftarrow C_{\Gamma_j^i}^{seg}(\mathbf{z}_j^i(\mathbf{p}) \odot \tilde{\mathbf{B}}_{j+1}^i(\mathbf{p}; 1))$
- 18:              $\mathbf{B}_j^i.append(\{\mathbf{B}_j^i(\mathbf{p}; \gamma)\}_{\gamma}^{\Gamma_j^i})$
- 19:         **else**  
*/\*Top-down masking for all other layers\*/*
- 20:             **for**  $\mathbf{B}_{j+1}^i(\mathbf{p}; k) \in \{\mathbf{B}_{j+1}^i\}_{\gamma}^{\Gamma_{j+1}^i}$  **do**  
*/\*Bilinear interpolate mask to feature shape\*/*
- 21:                  $\tilde{\mathbf{B}}_{j+1}^i(\mathbf{p}; k) \leftarrow \text{BiInterp}(\mathbf{B}_{j+1}^i(\mathbf{p}; k), \mathbf{z}_j^i.shape)$   
*/\*Mask with higher layer binary mask, Eq. (2)\*/*
- 22:                  $\Gamma_j^i \leftarrow \text{silhouette}(\mathbf{z}_j^i(\mathbf{p}) \odot \tilde{\mathbf{B}}_{j+1}^i(\mathbf{p}; k))$
- 23:                  $\{\mathbf{B}_j^i(\mathbf{p}; \gamma)\}_{\gamma}^{\Gamma_j^i} \leftarrow C_{\Gamma_j^i}^{seg}(\mathbf{z}_j^i(\mathbf{p}) \odot \tilde{\mathbf{B}}_{j+1}^i(\mathbf{p}; k))$
- 24:                  $\mathbf{B}_j^i.append(\{\mathbf{B}_j^i(\mathbf{p}; \gamma)\}_{\gamma}^{\Gamma_j^i})$
- 25:             **end for**
- 26:         **end if**  
*/\*Create and save RGB Masks, Eq. (3)\*/*
- 27:         **for**  $\mathbf{B}_j^i(\mathbf{p}; \gamma) \in \mathbf{B}_j^i$  **do**
- 28:              $\mathbf{M}_j^i(\mathbf{p}; \gamma) \leftarrow \mathcal{I}^i \odot \mathbf{B}_j^i(\mathbf{p}; \gamma)$
- 29:              $\mathbf{M}_j.append(\mathbf{M}_j^i(\mathbf{p}; \gamma))$
- 30:         **end for**
- 31:          $\mathbf{M}.append(\mathbf{M}_j)$
- 32:     **end for**
- 33: **end for**  
**Return**  $\mathbf{M}$

---

	ResNet50			VGG16			MViT			ViT-b		
	RF	ACE	Ours	RF	ACE	Ours	RF	ACE	Ours	RF	ACE	Ours
Layer1	43	2.4	4.2	10	2.8	3.5	224	1.7	7.5	224	1.3	10.0
Layer2	99	2.0	11.9	32	2.4	6.8	224	1.0	15.2	224	1.7	24.4
Layer3	211	1.92	23.9	80	2.2	15.5	224	2.5	28.1	224	2.4	35.8
Layer4	435	2.1	47.9	176	2.2	46.8	224	2.4	50.1	224	2.4	54.1

Table 1. Validation of segment proposal component of our method (Sec. 3.1). The relative concept segment size compared to the entire image for a given layer, is shown with the receptive field (RF) width/height of the same layer. We compare our method (Ours) to the baseline method, ACE [23]. For all models, the relative segment size discovered using our method has a stronger correlation with the receptive field size than the concepts discovered using ACE.

---

### Algorithm 2 Concept Discovery

---

**Input:** Model  $F$ ,  $n$  selected layers of  $F$  to study, Set of RGB Image Masks at each Layer  $\mathbf{M} = \{\mathbf{M}_1, \dots, \mathbf{M}_n\}$ , Clustering algorithm  $C^{con}$

*/\* $C^{con}$  is instantiated in terms of k-means [35]\*/*

**Output:** Set of concept centroids  $\mathbf{Q} = \{\mathbf{Q}_1, \dots, \mathbf{Q}_n\}$

```

1:  $\mathbf{Q} \leftarrow []$ 
2: for  $j \in n$  do
   /*Cluster segment activations, Eq. (4)*/
3:    $\mathbf{Z}_{\mathbf{M}_j} \leftarrow f_j(\mathbf{M}_j)$ 
4:    $\mathbf{Q}_j \leftarrow C^{con}(\text{GAP}(\mathbf{Z}_{\mathbf{M}_j}))$ 
   /*Prune clusters (see Sec. 4.1) for details*/
5:    $\mathbf{Q}_j \leftarrow \text{prune}(\mathbf{Q}_j)$ 
6:    $\mathbf{Q}.\text{append}(\mathbf{q}_j)$ 
7: end for
Return  $\mathbf{Q}$ 

```

---

- ResNet18 [25], ResNet50 [25] and CLIP-ResNet50 [42]: *Layer1, Layer2, Layer3, Layer4*
- VGG16 [48]: 8, 15, 22, 29
- MobileNetv3 [44]: 0, 2, 4, 6
- MViT [18]: 1, 3, 9, 15
- ViT-b [51]: 2, 5, 8, 10

The layer names used (according to the PyTorch [41] module nomenclature) for each model when generating the all-layer VCCs are as follows:

- ResNet50 [42]: *layer1.0, layer1.1, layer1.2, layer2.0, layer2.1, layer2.2, layer2.3, layer3.0, layer3.1, layer3.2, layer3.3, layer3.4, layer3.5, layer4.0, layer4.1, layer4.2*
- VGG16 [48]: 1, 3, 6, 8, 11, 13, 15, 18, 20, 22, 25, 27, 29
- MobileNetv3 [44]: 0.0, 1.0, 1.1, 2.0, 2.1, 2.2, 3.0, 3.1, 3.2, 3.3, 4.0, 4.1, 5.0, 5.1, 5.2, 6.0
- MViT [18]: 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15
- ViT-b [51]: 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10

**Clustering details.** Following previous work [23], during the concept discovery clustering,  $C^{con}$ , we over-cluster and then prune to ensure that fewer concepts will be missed. We follow previous work [23] and choose the number of

---

### Algorithm 3 Interlayer Testing with Concept Activation Vectors

---

**Input:** Model  $F$ , higher layer selected to study  $l$ , lower layer selected to study  $j$ , Concept Centroid for higher layer  $\mathbf{q}_l^{m_l}$ , Concept centroid for lower layer  $\mathbf{q}_j^{m_j}$ , Set of RGB image masks each associated with higher layer concept centroid  $\mathbf{M}_{\mathbf{q}_l^{m_l}}$ , Set of RGB image masks each associated with lower layer concept centroid  $\mathbf{M}_{\mathbf{q}_j^{m_j}}$ , Set of random images  $\mathcal{I}_{rnd}$ , Linear classifier  $h$

**Output:** Concept connection edge weight between concepts  $\mathbf{q}_j^{m_j}$  and  $\mathbf{q}_l^{m_l}$ :  $e_{\mathbf{q}_j^{m_j}, \mathbf{q}_l^{m_l}}$

*/\*Get activations for lower level concept\*/*

```

1:  $\mathbf{z}_{\mathbf{M}_{\mathbf{q}_j^{m_j}}} \leftarrow f_j(\mathbf{M}_{\mathbf{q}_j^{m_j}})$ 
   /*Get activations for random concept*/
2:  $\mathbf{z}_{\mathcal{I}_{rnd}} \leftarrow f_j(\mathcal{I}_{rnd})$ 
   /*Train CAV and get orthogonal vector to hyperplane in direction of lower concept*/
3:  $\mathbf{V}_{\mathbf{q}_j^{m_j}} \leftarrow h(\mathbf{z}_{\mathbf{M}_{\mathbf{q}_j^{m_j}}}, \mathbf{z}_{\mathcal{I}_{rnd}}).\text{train}()$ 
4:  $\text{CountPositive} \leftarrow 0$ 
   /*Iterate through higher concept segments*/
5: for  $x \in \mathbf{M}_{\mathbf{q}_l^{m_l}}$  do
6:    $\mathbf{z}_j \leftarrow f_j(x)$ 
   /*Get gradient of segment at layer l with respect to lower layer j*/
7:    $\mathbf{g}_j \leftarrow \nabla_{f_j} \|f_l(\mathbf{z}_j) - \mathbf{q}_l^{m_l}\|_2$ 
   /*Calculate sensitivity of upper concept to lower concept, Eq. (5)*/
8:    $S_{\mathbf{q}_j^{m_j}, \mathbf{q}_l^{m_l}} = \mathbf{g}_j \cdot \mathbf{V}_{\mathbf{q}_j^{m_j}}$ 
9:   if  $S_{\mathbf{q}_j^{m_j}, \mathbf{q}_l^{m_l}} > 1$  then
10:      $\text{CountPositive} = \text{CountPositive} + 1$ 
11:   end if
12: end for
   /*Calculate fraction of positive alignments, Eq. (6)*/
13:  $e_{\mathbf{q}_j^{m_j}, \mathbf{q}_l^{m_l}} = \text{CountPositive} / |\mathbf{M}_{\mathbf{q}_l^{m_l}}|$ 
Return  $e_{\mathbf{q}_j^{m_j}, \mathbf{q}_l^{m_l}}$ 

```

---

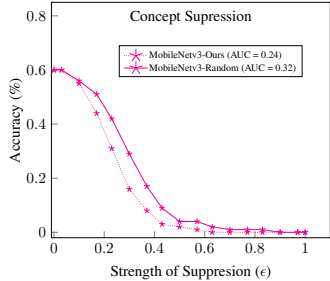


Figure 9. Additional validation results for the concept discovery method (Sec. 3.2 in the main paper) for the MobileNetV3 model [26]. For a set of 50 randomly selected ImageNet classes, we discover concepts in four layers of the model. During inference, one randomly selected concept at each layer is suppressed by a factor of  $\epsilon$ .

clusters to be  $k_m = 25$  in the concept discovery step. However, as VCCs target the discovery of concepts at potentially all layers, we select a different pruning protocol [23], where they prune based on a single minimum value. Instead, we prune clusters that have less than  $Y$  images, via the generalized logistic sigmoid

$$Y = A + \frac{K - A}{(C + Qe^{-Bt})^{1/\nu}}, \quad (9)$$

where  $A = -102$ ,  $K = 115$ ,  $C = 1$ ,  $Q = 1$ ,  $B = 0.0004$  and  $\nu = 1$ . Pruning based on a sigmoid shaped function enables different levels of leniency when considering what constitutes a concept for each layer. This is crucial as different layers contain a different number of segments from the top-down segmentation algorithm (Sec. 3.1 in the main paper).

For the maskSLIC [28] clustering stage, we use a compactness of 0.8 and all other parameters are set to the Scikit-Image [50] defaults. The Euclidean distance is used for all clustering steps.

**Randomized testing.** When applying our ITCAV method to calculate the strength of connection between two concepts, we protect against the impact of spurious results by performing a statistical significance test on all ITCAV scores. More specifically, instead of simply calculating the ITCAV score with the target concept images, we calculate an additional 20 ITCAV scores using random images sampled from Broden [3]. We perform a two-sided  $t$ -test of the ITCAV score based on the 20 random scores. We test whether the null hypothesis (*i.e.* a ITCAV score of 0.5) can be rejected with a  $p$ -value of  $p > 0.05$ . All ITCAV scores shown in the main paper and appendix pass this statistical test, *i.e.*  $p \leq 0.05$ .

**ImageNet classes.** The 50 ImageNet classes used for the model and task analysis experiments (Sec. 4.3 in the main paper) are the following: *tow truck, sturgeon,*

*sax, wool, basketball, whiptail, toy poodle, acorn, crutch, church, backpack, spaghetti squash, snowmobile, teapot, ant, chain, gorilla, holster, wreck, ice lolly schipperke, cradle, dowitcher, leopard, oystercatcher, saltshaker, drake, loupe, spotlight, Newfoundland, bagel, electric fan, ping-pong ball, streetcar, knot, plate, sea lion, leafhopper, tusker, punching bag, black widow, traffic light, tricycle, paper towel, guinea pig, castle, go-kart, platypus, badger and bicycle-built-for-two.*

The 10 ImageNet classes used for the all-layer VCC analysis experiments (Sec. 4.3 in the main paper) are the following: *tow truck, sturgeon, sax, wool, basketball, whiptail, toy poodle, acorn, crutch, church*

## 10. Additional empirical results

### 10.1. VCC component validation

#### 10.1.1 Segment proposal validation

Table 1 presents additional results to validate our top-down feature segmentation approach (Sec. 3.1). In particular, we show results for four additional models. We observe findings consistent with those in the main paper: Our method produces concepts that increase in size as the information flows deeper through the model. It is interesting to observe a similar phenomenon in transformer-based architectures, *i.e.* MViT [18] and ViT-b [51]. While the relative concept size of the baseline, ACE [23], varies less than 2% across all architectures and layers, the size of concepts produced by our method can differ up to 20% between architectures (*i.e.* comparing VGG16-Ours with ViT-b-Ours) and 40% between layers. This finding is to be expected as it is unlikely for all architectures at all layers to capture concepts of the same size.

#### 10.1.2 Concept validation

Figure 9 presents additional results to validate our layer-wise concept discovery method (Sec. 3.2 in the main paper) for the MobileNetV3 [26] model. The results are consistent with those in the main paper, *i.e.* the accuracy for the target class decreases faster when a concept is suppressed compared to a randomly chosen direction. This result implies that the concepts discovered throughout the model represent meaningful directions in the latent space.

#### 10.1.3 Interlayer concept weight validation

We now extend the validation of our Interlayer Testing with Concept Activation Vectors (ITCAV) method (Sec. 3.3 in the main paper). In particular, we show results for four additional models in Fig. 10 and observe findings consistent with those in the main paper: There is a positive correlation between the average path strength (APS) and the logit sum

	Branching Factor		Number of Concepts		Edge Weight Ave.	
	R50	CLIP	R50	CLIP	R50	CLIP
Layer1	5.484	6.824	10.447	11.085	0.414	0.417
Layer2	4.945	4.141	8.000	7.468	0.554	0.54
Layer3	2.799	2.754	5.106	5.702	0.476	0.563
Layer4	2.915	1.574	2.957	2.383	0.917	0.634

Table 2. VCC metrics for ResNet50 [25] trained on ImageNet [15] and via contrastive image-language pretraining (CLIP) [42].

(LS) score. These results further suggest that the combination of ITCAV scores is predictive of whether a concept is representative of the target class.

## 10.2. Understanding models

Figure 11 extends the results from Fig. 6 in the main paper and shows a quantitative analysis for all-layer VCCs on two additional models: ResNet50 [25] and ViT-b [16]. Consistent with the results from the main paper, we again see that the branching pattern and number of concepts start at a higher value and converges, suggesting that many concepts are shared between classes at early layers while the later layers capture ImageNet’s foreground-background structure. We also observe patterns in the ITCAV values and variances that are consistent with the main paper. The edge weight values are consistent until the final layer at which point they increase, denoting the stronger contribution of the final layers to the output. In terms of the ITCAV variance, we again see that transformers (ViT-b) have a higher variance than CNNs (ResNet50) in the last layer, further suggesting that transformers have greater compositionality of concepts before the final prediction.

## 10.3. Understanding tasks

We now explore how VCCs can reveal the effect of the training task on learned concepts and their connections. In particular, given the recent advances of image-language training paradigms, we compare the standard ResNet50 [25] model trained on ImageNet [15] with ResNet50 trained via Contrastive Language Image Pretraining (CLIP) [42].

Table 2 compares graph metrics over VCC layers between the two models at four residual blocks. We observe small but notable differences between the two models. First, CLIP contains a higher branching factor and number of concepts in the first layer than ResNet50, suggesting slightly more concepts are discovered and composed at the beginning of the network. The pattern is reversed at the end of the models, where CLIP has a slightly lower number of concepts and branching factor than ResNet50. When considering average edge weight values, we also observe a general consistency across models apart from the final layer, where ResNet50 has a much larger average value. This may be due to ImageNet trained CNNs having less compositionality at the end of the model as we observed both object and background classes having a large impact on the output in the main paper (Sec. 4.3).

## 10.4. Additional VCCs visualizations

### 10.4.1 Four layer VCCs

We now supplement the analysis from Sec. 4.3 from the main paper by generating VCCs for the entirety of the five models analyzed for different classes in the four layer setting. We specifically chose these models and layer settings as they are the same as in Sec. 4.3 in the main paper. The models shown are ResNet50 [25] (Fig. 13), VGG16 [48] (Fig. 14), MobileNetv3 [26] (Fig. 15), MViT [18] (Fig. 16) and ViT-b [51] (Fig. 17). All models are trained on ImageNet [15]. The layers selected are the same ones as detailed in Sec. 9.

We observe differences in mid and late layer connection strengths between CNNs and transformers. Similar to the main paper discussion (Sec. 4.3), CNNs (Figs. 13, 14 and 15) show stronger connections with less variance between the 4<sup>th</sup> layer and class logit than the transformers (Figs. 16 and 17). Additionally, CNNs tend toward concepts which capture either the entire foreground or background in later layers. Meanwhile, the transformers produce concept shapes of varying shapes and sizes, *e.g.* the VCC for ViT-b in Fig. 17 contains concepts of both small patches and the entire images in the final VCC layer and concepts of varying sizes in the first VCC layer. These findings for the transformers are consistent with the ability of such models to form data associations across their input without the locality constraints that are inherent in convolutional models.

**Finegrained dataset VCC.** To show how VCCs generalize to other datasets, we generate a VCC for the CUB [52] finegrained classification dataset, where the goal is to classify different types of birds. Figure 18 shows a four layer VCC for the ResNet18 [25] model targeting the class “indigo bunting”. We again see interesting concepts being composed. For example, branches and the color blue occur in stage1 and stage2, while stage4 bird concepts are composed from branch, background and bird head concepts in stage3.

**All layer VCC.** We show an additional all-layer VCC in Fig. 19 of the VGG16 model [49] targeting recognition of class “church”. As in the visualization in the main paper, we visualize VCC subgraphs and observe interesting compositions occurring at different levels of abstraction corresponding at different depths of the model. At early layers (bottom left), we observe oriented brown patterns and yellow color composing the concept of brown and yellow orientation. Middle layers (right) show the concept of ‘church roof with sky in the background’ being composed of ‘church roof’ and ‘sky’. The final layer concepts (top left) show that both foreground objects, *e.g.* churches, and background regions, *e.g.* trees or sky, concepts highly influence the final category.

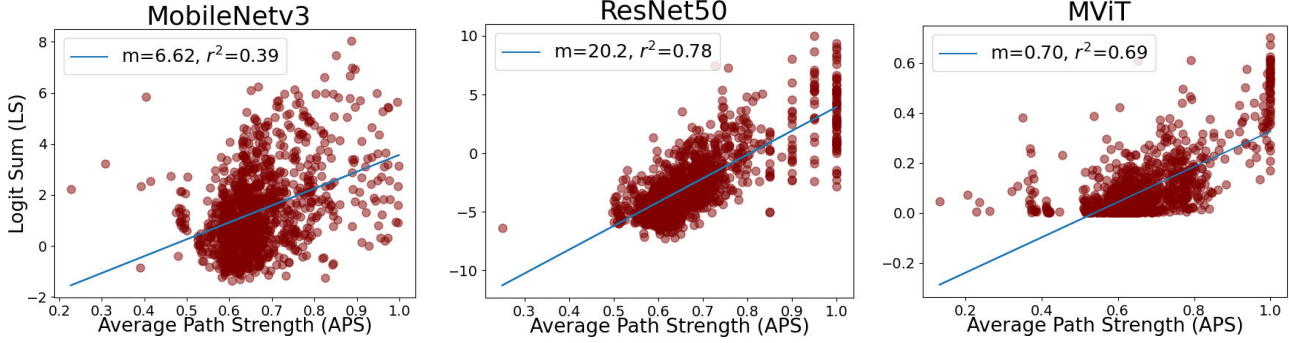


Figure 10. Additional validation results of interlayer concept weights. The unnormalized logit sum (LS) scores, main paper (Eq. 8), for the target class are plotted against the average path strength (APS) scores, main paper (Eq. 7). A positive correlation implies that the ITCAV edge weights connecting a concept to the class are predictive of the model output having a higher probability for that class.

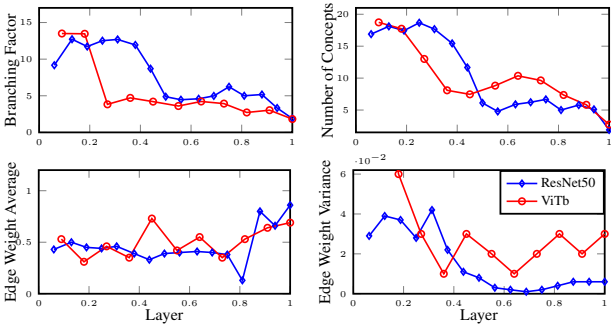


Figure 11. Graph metrics of all layer VCCs for two architectures. Layer number normalized to allow for comparison of models with different numbers of layers.

## 11. Application: Diagnosing failure predictions

To further show our VCC’s practical utility, we show another example of model debugging. Figure 12 shows a ‘church’ incorrectly predicted as a ‘vault’ by a ResNet50 model, and the corresponding incorrect VCC (‘vault’, left) and correct VCC (‘church’, right). As the image is decomposed using our top-down segmentation (Sec. 3.1), it is revealed that several segments are closer, in terms of the  $l_2$  distance of the pooled segment activations, to concepts in the ‘vault’ VCC (red outlines) than the ‘church’ VCC (green outlines). While the model correctly encoded the door as a ‘church’ concept, the regions outside the door are identified as ‘vault’ concepts from layers two through four, which may cause the error. We also note the lack of other ‘church’ specific concepts, such as the sky or cylindrical columns.

## 12. Limitations

We note some limitations of our method. We rely on the Silhouette method [43] to select the number of clusters (*i.e.* segments) during the top-down feature segmentation stage to automate this step. However, use of a different method for selecting the number of clusters could yield different

results and therefore different overall VCCs. In practice, we have found that using the Silhouette method consistently produces meaningful segments; so, sensitivity to this choice is not a serious limitation. Another limitation arises is that we do not provide a method for selecting the set of layers to analyze. Such a method for automatic layer selection could reveal further interesting and useful patterns, such as uncovering the set of layers, along with their connections, which impact the model output most significantly. A direction to realize such an algorithm could be to construct a large VCC and subsequently trim the least important nodes and edges (*e.g.* based on the average path strength (APS) to the logit, as defined in Eq. 7 in the main paper).

## 13. Societal implications

Understanding the decision making processes of deep networks is an important and open problem in computer vision. Given their potential for negative impacts when deployed, various jurisdictions are moving forward with legislation that may curtail certain applications and mandate interpretable components in deployed systems [14]. VCCs are a step towards a holistic understanding of how concepts in deep networks are learned and in the future may provide a direction to design legally recognized interpretations of these models.

VCCs may have implications in terms of recognizing both *what* and *how* biases are learned by deep networks. While the learning of various biases by deep networks is well documented [37], it is not well understood *how* these biases are constructed and learned by the model. For example, it is not sufficient to explain a model’s prediction by saying it uses the background as a feature. It would be more desirable to explain what concepts are composed in earlier layers that lead the model to encode the background feature in the later layers, which we have shown that VCCs can reveal. Moreover, such information could open up new directions for model debiasing.

In terms of negative consequences, VCCs (and explain-

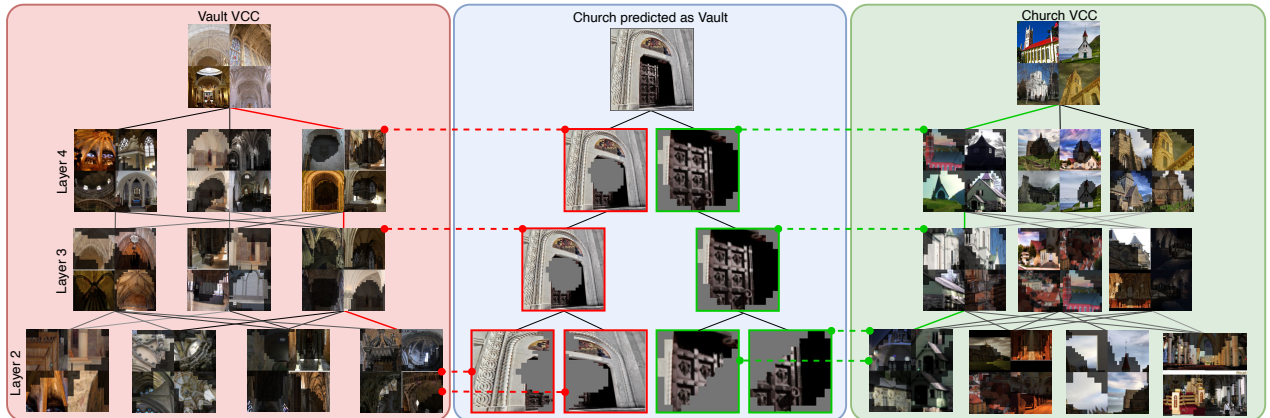


Figure 12. Debugging model failure modes with VCCs. We show an image of a church incorrectly predicted by a ResNet50 as a vault (middle) as well as the top-down segmentation of the image (Sec. 3.1). We also show the incorrect (left) and correct (right) VCCs. Following the hierarchy of concepts reveals that the model incorrectly focused heavily on the cement door frame, starting at Layer 2.

able AI in general) may give users a false sense of security and allow them to deploy models that ultimately do more harm than good. Furthermore, the contribution of additional explainable AI methods may contribute to the disagreement problem [33], *i.e.* when multiple explanations of a given model disagree with each other. It is an open research question on how to resolve such disagreements, when potentially dozens of possible explanations for a given model exist.

## 14. Assets and licensing

**Models.** We use provided code and trained weights from the MViT<sup>1</sup> and CLIP<sup>2</sup> repositories. MViT is licensed under the Apache 2.0 license<sup>3</sup> and CLIP is licensed under the MIT license<sup>4</sup>.

**Datasets.** We use the ImageNet dataset<sup>5</sup> which is under the BSD 3-Clause License<sup>6</sup> and the Broden dataset<sup>7</sup> which is under the MIT license<sup>8</sup>.

<sup>1</sup><https://github.com/facebookresearch/mvit>

<sup>2</sup><https://github.com/openai/CLIP>

<sup>3</sup><https://github.com/facebookresearch/mvit/blob/main/LICENSE>

<sup>4</sup><https://github.com/openai/CLIP/blob/main/LICENSE>

<sup>5</sup><https://www.image-net.org/>

<sup>6</sup><https://github.com/floydhub/imagenet/blob/master/LICENSE>

<sup>7</sup><https://github.com/CSAILVision/NetDissect-Lite>

<sup>8</sup><https://github.com/davidbau/quick-netdissect/blob/master/LICENSE>

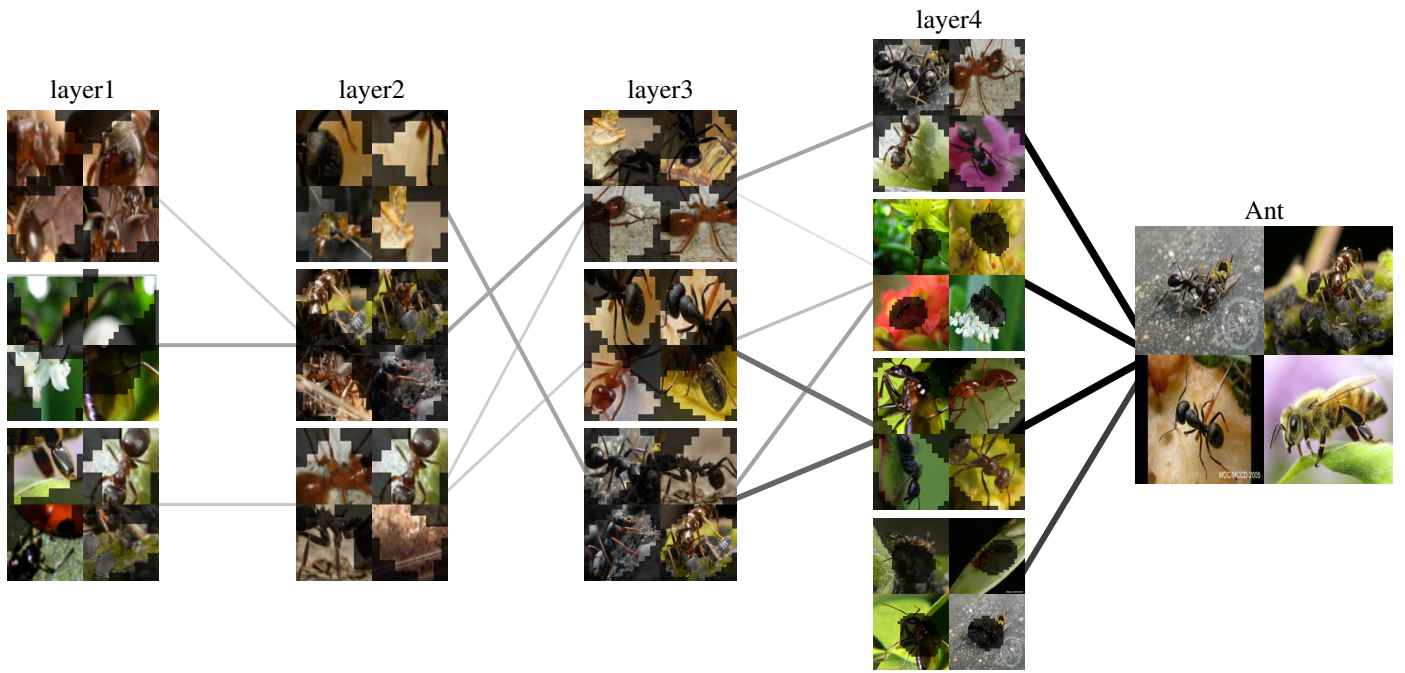


Figure 13. A VCC for four layers of a ResNet50 [25] model targeting the class “ant”. Darker lines denote larger concept contributions.

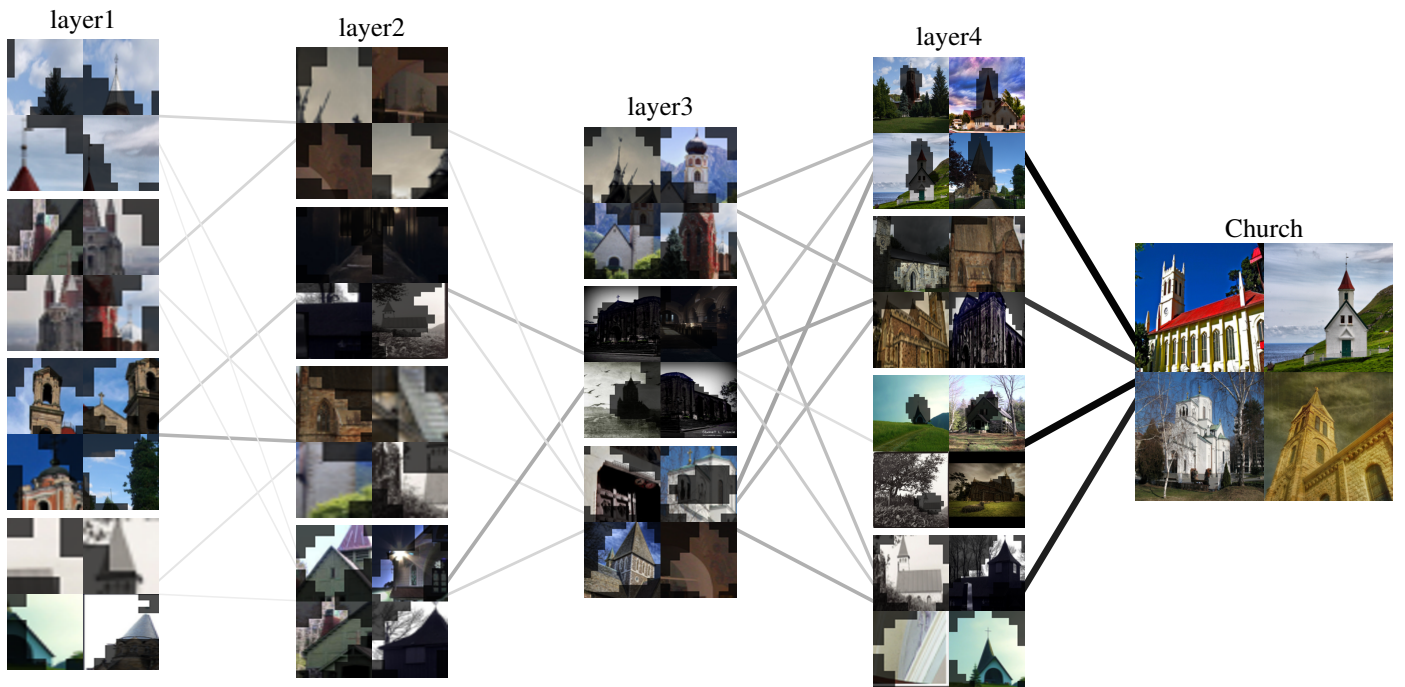


Figure 14. A VCC for four layers of a VGG16 [49] model targeting the class “church”. Darker lines denote larger concept contributions.

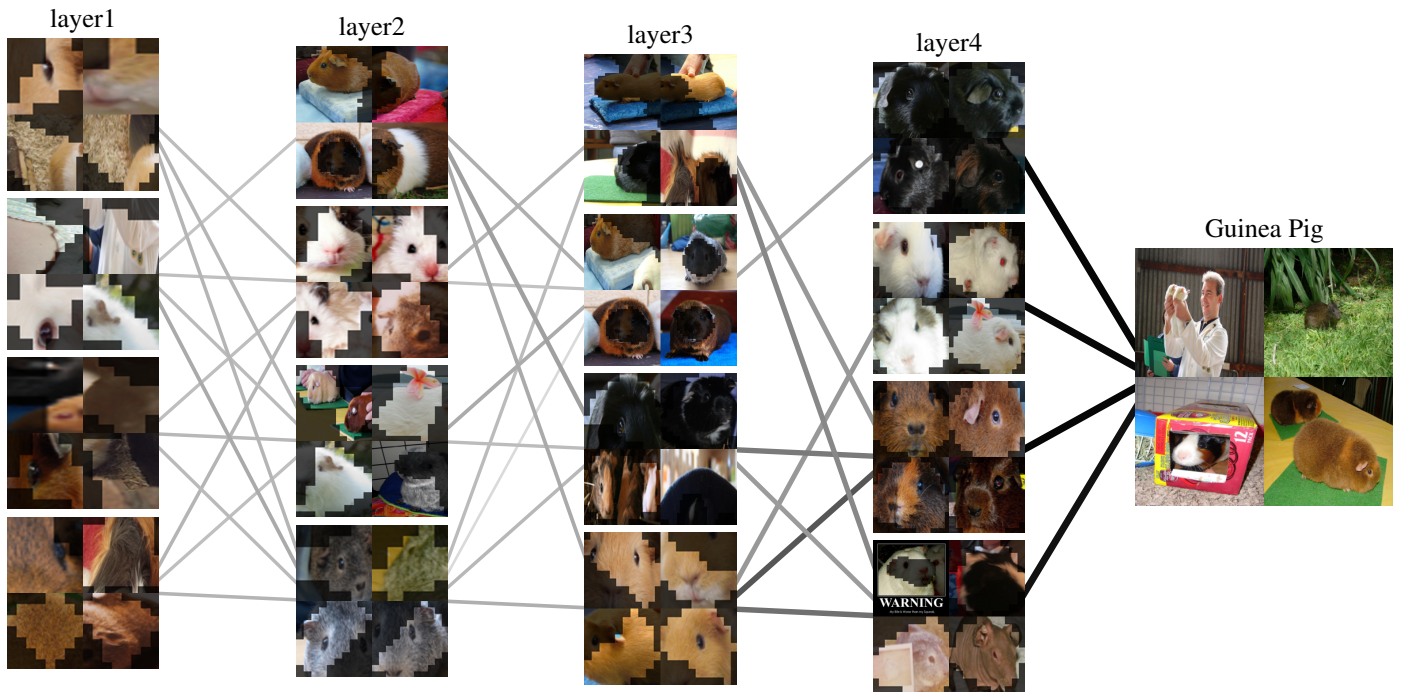


Figure 15. A VCC for four layers of a MobileNetv3 [26] model targeting the class “guinea pig”. Darker lines denote larger concept contributions.

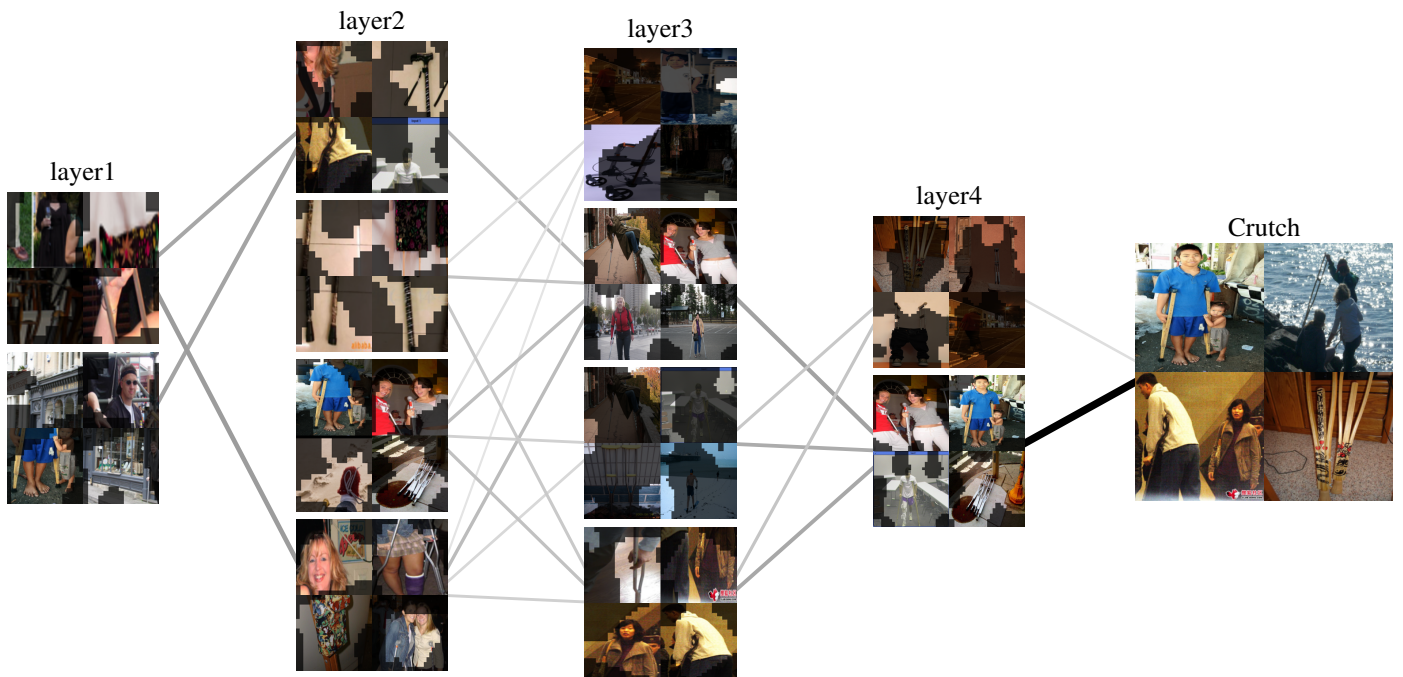


Figure 16. A VCC for four layers of a MViT [18] model targeting the class “crutch”. Darker lines denote larger concept contributions.



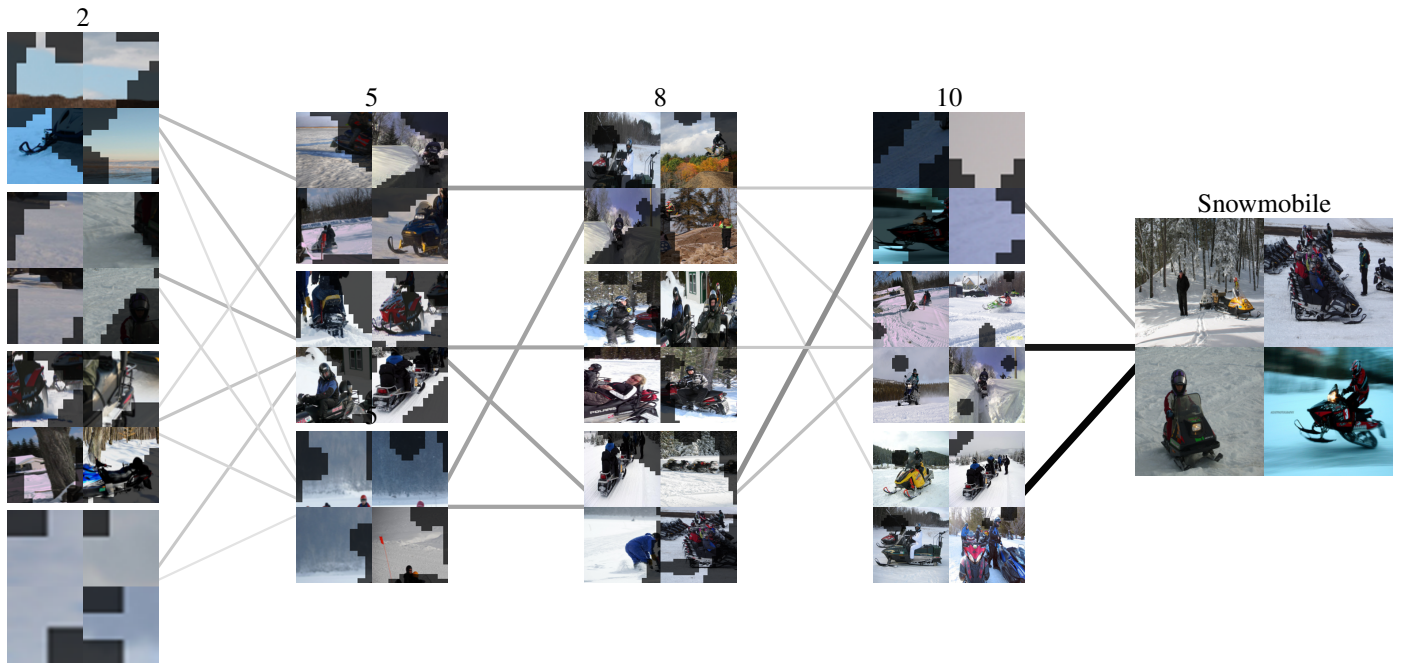


Figure 17. A VCC for four layers of a ViT [16] model targeting the class “snowmobile”. Darker lines denote larger concept contributions.

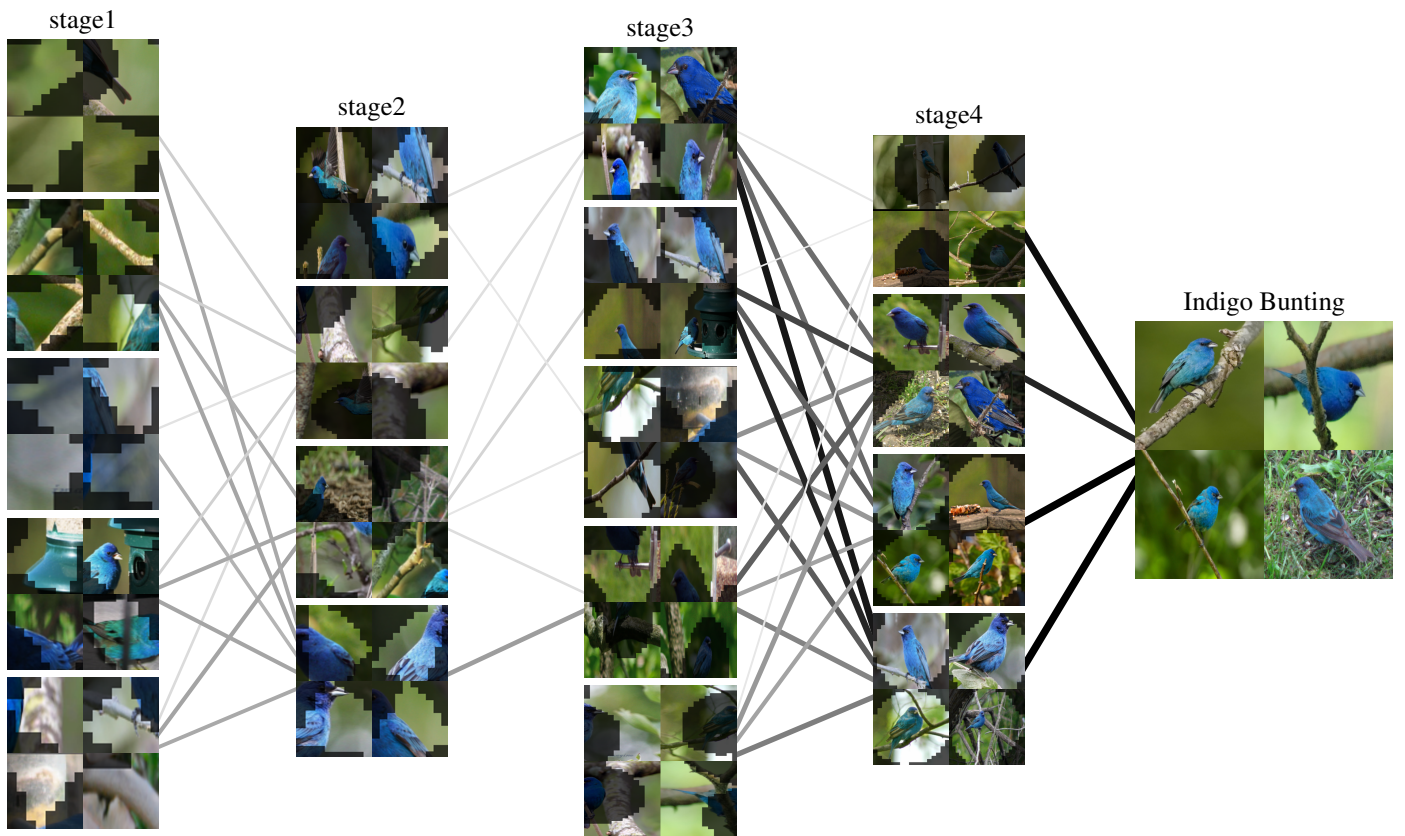


Figure 18. A VCC for four layers of a ResNet18 [25] model trained on the finegrained CUB [52] dataset, targeting the class “indigo bunting”. Darker lines denote larger concept contributions.

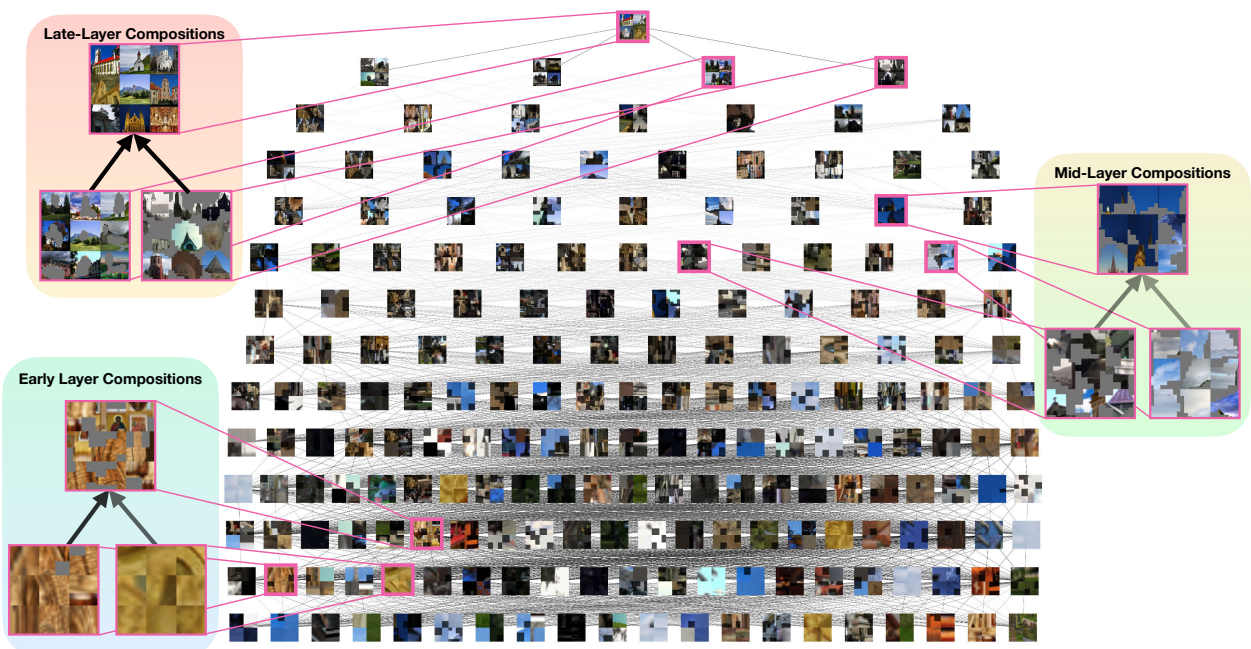


Figure 19. An all-layer VCC of the VGG16 network targeting the class “church”. Darker lines denote larger concept contributions.

## References

- [1] Reduan Achtibat, Maximilian Dreyer, Ilona Eisenbraun, Sebastian Bosse, Thomas Wiegand, Wojciech Samek, and Sebastian Lapuschkin. From “where” to “what”: Towards human-understandable explanations through concept relevance propagation. *arXiv preprint arXiv:2206.03208*, 2022. **3**
- [2] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*, 2018. **2**
- [3] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Conference on Computer Vision and Pattern Recognition*, pages 6541–6549, 2017. **2, 5, 1, 3**
- [4] David Bau, Steven Liu, Tongzhou Wang, Jun-Yan Zhu, and Antonio Torralba. Rewriting a deep generative model. In *European Conference on Computer Vision*, pages 351–369. Springer, 2020. **2**
- [5] Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, Klaus-Robert Müller, and Wojciech Samek. Layer-wise relevance propagation for neural networks with local renormalization layers. In *Artificial Neural Networks and Machine Learning*, pages 63–71. Springer, 2016. **2**
- [6] Moritz Böhle, Mario Fritz, and Bernt Schiele. B-cos networks: Alignment is all we need for interpretability. In *Conference on Computer Vision and Pattern Recognition*, pages 10329–10338, 2022. **2**
- [7] Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermy, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023. <https://transformer-circuits.pub/2023/monosemantic-features/index.html>. **3**
- [8] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on Fairness, Accountability and Transparency*, pages 77–91, 2018. **1**
- [9] Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-CAM++: Generalized gradient-based visual explanations for deep convolutional networks. In *Winter Conference on Applications of Computer Vision*, pages 839–847, 2018. **2**
- [10] Hila Chefer, Shir Gur, and Lior Wolf. Transformer interpretability beyond attention visualization. In *Conference on Computer Vision and Pattern Recognition*, pages 782–791, 2021. **1**
- [11] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. **7**
- [12] Zhi Chen, Yijie Bei, and Cynthia Rudin. Concept whitening for interpretable image recognition. *Nature Machine Intelligence*, 2(12):772–782, 2020. **2**
- [13] Ming-Ming Cheng, Peng-Tao Jiang, Ling-Hao Han, Liang Wang, and Philip Torr. Deeply explain CNN via hierarchical decomposition. *International Journal of Computer Vision*, 131:1091–1105, 2023. **3, 8**
- [14] European Commission. Laying down harmonised rules on artificial intelligence (artificial intelligence act) and amending certain union legislative acts. *European Commission*, 2021. **5**
- [15] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. **2, 6, 1, 4**
- [16] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, and Sylvain Gelly. An image is worth 16x16 words: Transformers for image recognition at scale. 2021. **6, 4, 9**
- [17] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition. *Transformer Circuits Thread*, 2022. [https://transformer-circuits.pub/2022/toy\\_model/index.html](https://transformer-circuits.pub/2022/toy_model/index.html). **3**
- [18] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. *International Conference on Computer Vision*, 2021. **5, 2, 3, 4, 8**
- [19] Christoph Feichtenhofer, Axel Pinz, Richard P Wildes, and Andrew Zisserman. Deep insights into convolutional networks for video recognition. *International Journal of Computer Vision*, 128:420–437, 2020. **1, 2**
- [20] Thomas Fel, Agustin Picard, Louis Bethune, Thibaut Boissin, David Vigouroux, Julien Colin, Rémi Cadène, and Thomas Serre. CRAFT: Concept recursive activation factorization for explainability. In *Conference on Computer Vision and Pattern Recognition*, pages 2711–2721, 2023. **1, 2, 3, 6**
- [21] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *International Conference on Learning Representations*, 2018. **7**
- [22] Amirata Ghorbani, Abubakar Abid, and James Zou. Interpretation of neural networks is fragile. In *AAAI Conference on Artificial Intelligence*, pages 3681–3688, 2019. **2**
- [23] Amirata Ghorbani, James Wexler, James Y Zou, and Been Kim. Towards automatic concept-based explanations. In *Advances in Neural Information Processing Systems*, 2019. **1, 2, 3, 4, 5, 6**
- [24] Sven Ove Hansson, Matts-Åke Belin, and Björn Lundgren. Self-driving vehicles-An ethical overview. *Philosophy & Technology*, 34:1383–1408, 2021. **1**
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference*

- on *Computer Vision and Pattern Recognition*, pages 770–778, 2016. 5, 6, 7, 2, 4, 9
- [26] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, and Vijay Vasudevan. Searching for MobileNetV3. In *International Conference on Computer Vision*, pages 1314–1324, 2019. 5, 3, 4, 8
- [27] Haiyang Huang, Zhi Chen, and Cynthia Rudin. SegDiscover: Visual concept discovery via unsupervised semantic segmentation. *arXiv preprint arXiv:2204.10926*, 2022. 2
- [28] Benjamin Irving. MaskSLIC: Regional superpixel generation with application to local pathology characterisation in medical images. *arXiv preprint arXiv:1606.09518*, 2016. 4, 1, 3
- [29] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, and Fernanda Viegas. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV). In *International Conference on Machine Learning*, pages 2668–2677, 2018. 1, 2, 4, 5, 6
- [30] Pieter-Jan Kindermans, Sara Hooker, Julius Adebayo, Maximilian Alber, Kristof T Schütt, Sven Dähne, Dumitru Erhan, and Been Kim. The (un) reliability of saliency methods. *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, pages 267–280, 2019. 2
- [31] Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. Concept bottleneck models. In *International Conference on Machine Learning*, pages 5338–5348, 2020. 2
- [32] Matthew Kowal, Achal Dave, Rares Ambrus, Adrien Gaidon, Konstantinos G Derpanis, and Pavel Tokmakov. Understanding video transformers via universal concept discovery. *arXiv preprint arXiv:2401.10831*, 2024. 1
- [33] Satyapriya Krishna, Tessa Han, Alex Gu, Javin Pombra, Shahin Jabbari, Steven Wu, and Himabindu Lakkaraju. The disagreement problem in explainable machine learning: A practitioner’s perspective. *arXiv preprint arXiv:2202.01602*, 2022. 6
- [34] Mengchen Liu, Jiaxin Shi, Kelei Cao, Jun Zhu, and Shixia Liu. Analyzing the training processes of deep generative models. *Transactions on Visualization and Computer Graphics*, 24(1):77–87, 2017. 2
- [35] Stuart Lloyd. Least squares quantization in PCM. *Transactions on information theory*, 28(2):129–137, 1982. 2
- [36] Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting them. In *Conference on Computer Vision and Pattern Recognition*, pages 5188–5196, 2015. 2
- [37] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)*, 54(6):1–35, 2021. 5
- [38] Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. *Distill*, 2017. <https://distill.pub/2017/feature-visualization>. 1, 2
- [39] Chris Olah, Arvind Satyanarayan, Ian Johnson, Shan Carter, Ludwig Schubert, Katherine Ye, and Alexander Mordvintsev. The building blocks of interpretability. *Distill*, 2018. <https://distill.pub/2018/building-blocks>. 1
- [40] Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. Zoom in: An introduction to circuits. *Distill*, 2020. <https://distill.pub/2020/circuits/zoom-in>. 2, 3
- [41] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, and Luca Antiga. PyTorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems*, 32, 2019. 1, 2
- [42] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763, 2021. 1, 2, 4
- [43] Peter J Rousseeuw. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20:53–65, 1987. 4, 5
- [44] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. MobileNetV2: Inverted residuals and linear bottlenecks. In *Conference on Computer Vision and Pattern Recognition*, pages 4510–4520, 2018. 2
- [45] Anirban Sarkar, Deepak Vijaykeerthy, Anindya Sarkar, and Vineeth N Balasubramanian. A framework for learning ante-hoc explainable models via concepts. In *Conference on Computer Vision and Pattern Recognition*, pages 10286–10295, 2022. 2
- [46] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-CAM: Visual explanations from deep networks via gradient-based localization. In *International Conference on Computer Vision*, pages 618–626, 2017. 1, 2, 8
- [47] S. Seung. *Connectome: How the Brain’s Wiring Makes Us Who We Are*. Houghton Mifflin Harcourt, 2012. 1
- [48] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *International Conference on Learning Representations*, 2014. 5, 6, 2, 4
- [49] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Conference on Computer Vision and Pattern Recognition*, pages 1–9, 2015. 2, 6, 7, 4
- [50] Stefan Van der Walt, Johannes L Schönberger, Juan Nunez-Iglesias, François Boulogne, Joshua D Warner, Neil Yager, Emmanuelle Gouillart, and Tony Yu. Scikit-image: image processing in python. *PeerJ*, 2:e453, 2014. 3
- [51] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017. 5, 7, 2, 3, 4
- [52] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech UCSD birds 2011 dataset. Technical

Report CNS-TR-2011-001, California Institute of Technology, 2011. [5](#), [4](#), [9](#)

- [53] Chih-Kuan Yeh, Been Kim, Sercan Arik, Chun-Liang Li, Tomas Pfister, and Pradeep Ravikumar. On completeness-aware concept-based explanations in deep neural networks. In *Advances in Neural Information Processing Systems*, pages 20554–20565, 2020. [2](#)
- [54] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision*, pages 818–833, 2014. [1](#), [2](#)
- [55] Quanshi Zhang, Ruiming Cao, Feng Shi, Ying Nian Wu, and Song-Chun Zhu. Interpreting CNN knowledge via an explanatory graph. In *AAAI Conference on Artificial Intelligence*, pages 4454–4463, 2018. [2](#)
- [56] Quanshi Zhang, Ying Nian Wu, and Song-Chun Zhu. Interpretable convolutional neural networks. In *Conference on Computer Vision and Pattern Recognition*, pages 8827–8836, 2018. [2](#)
- [57] Quanshi Zhang, Yu Yang, Haotian Ma, and Ying Nian Wu. Interpreting CNNs via decision trees. In *Conference on Computer Vision and Pattern Recognition*, pages 6261–6270, 2019. [2](#)
- [58] Ruihan Zhang, Prashan Madumal, Tim Miller, Krista A Ehinger, and Benjamin IP Rubinstein. Invertible concept-based explanations for CNN models with non-negative concept activation vectors. In *AAAI Conference on Artificial Intelligence*, pages 11682–11690, 2021. [2](#)
- [59] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Conference on Computer Vision and Pattern Recognition*, pages 2921–2929, 2016. [1](#), [2](#)
- [60] Bolei Zhou, Yiyou Sun, David Bau, and Antonio Torralba. Interpretable basis decomposition for visual explanation. In *European Conference on Computer Vision*, pages 119–134, 2018. [2](#)