
Visual-TCAV: Explainability of Image Classification through Concept-based Saliency Maps

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Convolutional Neural Networks (CNNs) have seen significant performance im-
2 provements in recent years. However, due to their size and complexity, their
3 decision-making process remains a black-box, leading to opacity and trust issues.
4 State-of-the-art saliency methods can generate local explanations that highlight the
5 area in the input image where a class is identified but do not explain how different
6 features contribute to the prediction. On the other hand, concept-based methods,
7 such as TCAV (Testing with Concept Activation Vectors), provide global explain-
8 ability, but cannot compute the attribution of a concept in a specific prediction nor
9 show the locations where the network detects these concepts. This paper introduces
10 a novel explainability framework, Visual-TCAV, which aims to bridge the gap
11 between these methods. Visual-TCAV uses Concept Activation Vectors (CAVs) to
12 generate saliency maps that show where concepts are recognized by the network.
13 Moreover, it can estimate the attribution of these concepts to the output of any
14 class using a generalization of Integrated Gradients. Visual-TCAV can provide
15 both local and global explanations for any CNN-based image classification model
16 without requiring any modifications. This framework is evaluated on widely used
17 CNNs and its validity is further confirmed through experiments where a ground
18 truth for explanations is known.

19 1 Introduction

20 Recent advancements in Deep Neural Networks (DNNs) have revolutionized the field of Artificial
21 Intelligence, and Convolutional Neural Networks (CNNs) have emerged as the state-of-the-art for
22 image classification due to their ability to learn complex patterns and features within images. However,
23 as the performance of these models has grown significantly over recent years, their complexity has
24 also increased. Consequently, it became a challenge to understand how these models produce their
25 classifications. This led to the widespread use of the term *black-box* to describe these models, as only
26 their inputs and outputs are known, while their internal mechanisms remain too complex for humans
27 to comprehend. The black-box problem results in a lack of transparency [29], which can undermine
28 trust in AI-based systems [12]. Indeed, blindly trusting AI poses serious ethical dilemmas, especially
29 in critical fields such as healthcare or autonomous driving in which image classification systems are
30 becoming increasingly employed [28, 3]. Additionally, debugging black-box models and identifying
31 biases becomes difficult without comprehending the process they use to make predictions. To this
32 end, the field of Explainable Artificial Intelligence (XAI) has made significant progress in developing
33 techniques for producing explanations of AI decisions. However, comprehending the specific features
34 or patterns that networks identify in an image and their precise impact on the prediction remains a
35 challenge. State-of-the-art approaches for local explainability (i.e., for individual predictions) use
36 saliency maps to locate where a class is identified in an input image, but they can't explain which
37 features led the model to its prediction. For instance, when analyzing an image of a golf ball, these

38 saliency methods cannot determine whether the golf ball was recognized by the spherical shape, the
39 dimples, or some other feature. Striving to cover this need, Kim et al. [11] introduced TCAV (Testing
40 with Concept Activation Vectors), a concept-based method that can discern whether a user-defined
41 concept (e.g., dimples, spherical) correlates positively with the output of a selected class. However,
42 TCAV is designed exclusively for global explainability (i.e., for explaining the general behavior of a
43 model) and therefore cannot measure the influence of a concept in a specific prediction or show the
44 locations within the input images where the networks recognize these concepts.

45 In this article, we introduce a novel explainability framework, Visual-TCAV, which integrates the core
46 principles of both saliency methods and concept-based approaches while aiming to overcome their
47 respective limitations. Visual-TCAV can be applied to any layer of a CNN model whose output is a
48 set of feature maps. Its main contributions are: (a) it provides visual explanations that show where
49 the network identifies user-defined concepts; (b) it can estimate the importance of these concepts to
50 the output of a selected class; (c) it can be used for both local and global explainability.

51 **2 Related Works**

52 In recent years, there has been a significant increase in the body of work exploring the explainability
53 of black-box models. For CNN-based image classification, state-of-the-art methods primarily focus
54 on providing explanations via saliency maps. These heatmaps highlight the most important regions
55 of the input image and therefore can be used to gain insights into how a model makes its decisions.
56 One approach for generating such visualizations involves studying the input-output relationship
57 of the model by creating a set of perturbed versions of the input and analyzing how the output
58 changes with each perturbation. Notable contributions to this approach include Local Interpretable
59 Model-Agnostic Explanations (LIME) [17], which uses random perturbations, and SHapley Additive
60 exPlanations (SHAP) [14], which estimates the importance of each pixel using Shapley values. A
61 different approach that instead tries to access the internal workings of the model was originally
62 proposed by Simonyan et al. [22] and consists of generating saliency maps based on the gradients
63 of the model output w.r.t. the input images. This idea led many researchers [24, 23] to investigate
64 how to exploit gradients to produce more accurate saliency maps. Selvaraju et al. [20] proposed a
65 method named Gradient-weighted Class Activation Mapping (Grad-CAM) that extracts the gradients
66 of the logits (i.e., raw pre-softmax predictions) w.r.t. the feature maps. It then uses a Global Average
67 Pooling (GAP) operation to transform these gradients into class-specific weights for each feature
68 map and performs a weighted sum of these feature maps to produce a class localization map, a
69 saliency map that highlights where a class is identified. Grad-CAM has gained considerable attention
70 and is extensively used for explaining convolutional networks. However, Sundararajan et al. [25]
71 demonstrated that gradients can saturate, leading to an inaccurate assessment of feature importance.
72 To address this issue, they introduced Integrated Gradients (IG), a method that calculates feature
73 attribution by integrating the gradients along a path from a baseline (e.g., a black image) to the
74 actual input image. Notable contributions of IG and its variants [10, 16, 30] include the ability to
75 provide fine-grained saliency maps (i.e., each pixel has its attribution) and adherence to the axiom of
76 completeness (i.e., the sum of the attributions of all pixels equals the logit value).

77 While saliency methods are effective and intuitive, they might not always provide a complete picture
78 of why a model made a certain decision. This is because these methods perform class localization,
79 but cannot explain which features led the model to recognize the highlighted class. Furthermore,
80 these techniques rely on per-pixel importance which can't be generalized across multiple instances, as
81 the position of these pixels is only meaningful for a specific input image. Consequently, they can only
82 explain one image at a time, preventing them from providing global explanations. To overcome these
83 limitations, Kim et al. [11] proposed Testing with Concept Activation Vectors (TCAV), a method that
84 investigates the correlations between user-defined concepts and the network's predictions using a set
85 of example images representing a concept. For instance, images of stripes can be used to determine
86 whether the network is sensitive to the "striped" concept for predicting the "zebra" class. This is
87 accomplished by calculating a Concept Activation Vector (CAV), which is a vector orthogonal to
88 the decision boundary of a linear classifier, typically Support Vector Machines (SVMs), trained to
89 differentiate between the feature maps of concept examples and random images. From this, a TCAV
90 score for any concept and model's layer can be computed using the signs of the dot products between
91 the CAV and the gradients of the loss w.r.t. the feature maps produced by images of a selected class.
92 TCAV is effective in detecting specific biases in neural networks (e.g., ethnicity-related) and can be

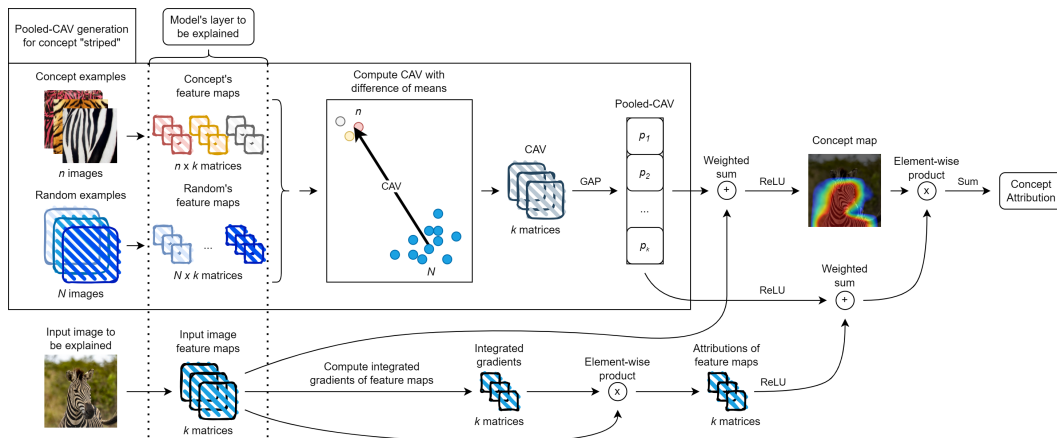


Figure 1: The Visual-TCAV process for generating local explanations. A Pooled-CAV is computed using the feature maps of user-defined concept examples and random images. A concept map is then produced through a weighted sum of the Pooled-CAV and the image’s feature maps. Finally, a concept attribution is obtained by extracting the IG attributions of the neurons that the concept activates using the Pooled-CAV and the concept map, which is used as a spatial mask.

93 considered complementary to saliency methods. Indeed, while saliency methods apply exclusively
 94 to individual predictions, TCAV can only provide global explanations. However, TCAV does not
 95 provide any information about the locations where concepts are identified within the input images.
 96 This makes it challenging to assess whether a high score can truly be attributed to the intended
 97 concept and not to a related one. Moreover, TCAV computes the network’s sensitivity to a concept,
 98 but not the magnitude of its importance in the prediction as the score only depends on the signs of the
 99 directional derivatives. For instance, “white” and “dimples” concepts might have identical TCAV
 100 scores for the “golf ball” class, even if one contributes substantially more to the prediction.

101 TCAV has received attention within the XAI community, leading to various extensions [5, 8] and
 102 applications [13, 2]. While our study focuses on user-defined concepts, unsupervised approaches
 103 have also been proposed. Ghorbani et al. [7] introduced Automatic Concept Extraction (ACE), a
 104 method that automatically extracts concepts from images for applying TCAV. This is accomplished
 105 by segmenting input images and subsequently clustering their activations. Building upon ACE, Zhang
 106 et al. [31] proposed Invertible Concept-based Explanations (ICE). This extension uses non-negative
 107 CAVs derived from non-negative matrix factorization and can also be used to explain locally by
 108 associating extracted concepts with a relevant area in the input image. Later, Bianchi et al. [1]
 109 proposed an unsupervised method for visualizing the entire feature extraction process of CNNs. They
 110 perform layer-wise clustering of similar feature maps to extract a set of concepts for each layer to
 111 which they assign a descriptive label through crowdsourcing. This approach provides local and global
 112 explanations, but the reliance on crowdsourcing can pose a practical challenge. Furthermore, these
 113 unsupervised approaches may provide opaque explanations. This is because, when the extracted
 114 image regions contain overlapping concepts (e.g., dimples, spherical, and white in a golf ball), it
 115 remains unclear which concepts the network has learned to recognize or considers more important.

116 3 Visual-TCAV

117 This section presents the methodology of our framework, Visual-TCAV, which is designed to explain
 118 the outputs of image classification CNNs using user-defined concepts. Local explanations can be
 119 generated considering any layer and consist of two key components. The first is the *Concept Map*, a
 120 saliency map that serves as a visual representation of the areas where the network has recognized
 121 the selected concept in the input image. The second is the *Concept Attribution*, a numerical value
 122 that estimates the importance of the concept for the output of a selected class. Figure 1 illustrates the
 123 pipeline for generating a local explanation. For global explanations, the process is replicated across
 124 multiple input images. The concept attributions for each image are then averaged to quantify how the
 125 concept influences the network’s decisions across a wide range of inputs.

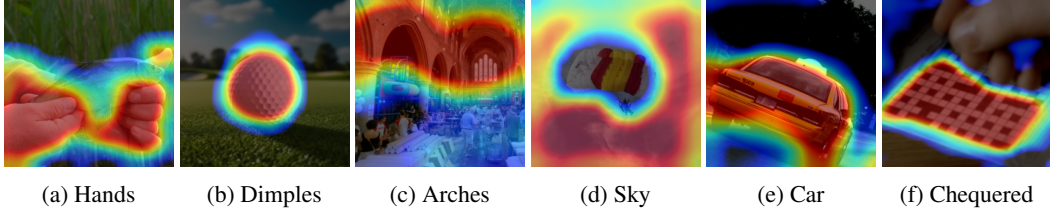


Figure 2: Examples of class-independent concept maps for various input images and concepts.

126 3.1 CAV Generation and Spatial Pooling

127 Similarly to the TCAV framework, the initial step of our method consists of computing a Concept
 128 Activation Vector (CAV) from a set of example images representing a user-defined concept, and a set
 129 of negative examples (e.g., random images). Specifically, we use the *Difference of Means* method,
 130 proposed by Martin and Weller [15], to compute the CAV. They demonstrated that this approach
 131 produces CAVs that are more resilient to perturbation and consistent than logistic classifiers or SVMs.
 132 As the name suggests, this method uses the arithmetic mean to determine the centroids of both the
 133 concept’s activations and the activations of random images. Subsequently, it directly computes the
 134 CAV as the difference between these centroids.

135 Since we are interested in identifying which feature maps are activated by the concept, irrespective of
 136 its location within the example images, we apply a Global Average Pooling (GAP) operation on the
 137 obtained CAV. The result is a vector of scalar values whose length is equal to the number of feature
 138 maps of the layer under consideration. Each vector element is associated with a feature map, and its
 139 raw value approximates the degree of correlation between that feature map and the concept. Moving
 140 forward, we will refer to this vector as the *Pooled-CAV*.

141 3.2 Concept Map

142 From the Pooled-CAV, we can construct a concept map that locates a concept (c) within any input
 143 image to be explained. This is achieved by performing a weighted sum of the feature maps ($fmaps_k$)
 144 of the input image, with the weights being the Pooled-CAV values (p_k^c). Equation (1) shows how
 145 to compute a raw concept map (M_{raw}^c). We also apply a ReLU function after the weighted sum
 146 because we are only interested in the image regions that positively correlate with the concept. The
 147 computation is similar to Grad-CAM’s equation, with the difference that we use the elements of the
 148 Pooled-CAV as weights instead of the global-average-pooled gradients.

$$M^{c,raw} = ReLU\left(\sum_k p_k^c \cdot fmaps_k\right) \quad (1)$$

149 We refer to this concept map as *raw* due to the absence of a scale factor (i.e., a maximum value) that
 150 would allow us to compare the degree of activation of the concept map across different concepts,
 151 input images, and model layers. To this end, we derive a concept map’s scale factor from the
 152 example images the user provided, which represent an ideal concept. Formally, we use Equation (2)
 153 to calculate the scale factor (s_c) as the maximum value of a hypothetical concept map, computed
 154 using the centroid (C^c), derived from the mean of the feature maps of the example images for a
 155 concept (c). Subsequently, we normalize the raw concept map by dividing it by the scale factor (s_c)
 156 and limiting the values to a unitary maximum, as shown in Equation (3). An epsilon (ϵ) is added to
 157 the denominator to prevent division by zero.

$$s_c = \max\left(ReLU\left(\sum_k p_k^c \cdot C_k^c\right)\right) \quad (2) \quad M_{ij}^c = \min\left(1, \frac{M_{ij}^{c,raw}}{s_c + \epsilon}\right) \quad \forall i, j \quad (3)$$

159 By overlaying the *normalized* concept map (M^c) on the input image, we can generate a class-
 160 independent visualization (examples are shown in Figure 2) that highlights the region of the image
 161 where the network recognized the concept. This allows us to know, for any input image, the
 162 concept’s location and its degree of activation w.r.t an ideal concept defined by the user. Additionally,

163 the concept map can provide a direct validation for the learned CAV, without requiring activation
 164 maximization techniques or sorting images based on their similarity to the CAV.

165 3.3 Concept Attribution

166 Once we acquire a set of concepts, we can gain insights into the network’s decision-making process
 167 by measuring the attribution of these user-defined concepts towards the raw predictions, also known
 168 as the logits. For instance, if the “church” class is predicted with a certain logit, we aim to quantify
 169 how much of this value is attributable to the “pews” concept, the “fresco” concept, and so on. More
 170 specifically, given an input image and a layer, we compute the attributions of the activations (i.e.,
 171 the values of the feature maps) to the logit of a specific target class. Subsequently, we utilize the
 172 Pooled-CAV to approximate which activations are attributable to a certain concept, and then we
 173 extract and sum these attributions. The attributions of a layer’s activations can be computed through a
 174 generalized variant of the IG approach which computes the integrated gradients of a target class’s
 175 logit w.r.t. the feature maps, instead of the input image. Specifically, we calculate the gradients along
 176 a straight-line path from zero-filled matrices to the actual feature maps and then approximate the
 177 integral using the Riemann trapezoidal rule. In our experiments, we consistently used 300 steps,
 178 which are sufficient to approximate the integral within a 5% error margin, as shown by Sundararajan
 179 et al. [25]. We then calculate the raw attributions by multiplying the integrated gradients with the
 180 feature maps, as shown in Figure 1. Since IG respects the completeness axiom regardless of which
 181 layer is considered as input, the attributions add up to the logit value of the target class, within the
 182 approximation error. A ReLU is then applied to extract positive attributions. These attributions
 183 are on the same scale as the raw logits, which can make their interpretation difficult. To obtain a
 184 comprehensible unitary scale, we normalize the attributions so that their sum equals a normalized
 185 logit, not the raw one. These normalized logits are obtained by applying a ReLU, followed by
 186 $[0,1]$ rescaling to retain their relative ratios.

187 To estimate the attribution of a concept (c), we can utilize the Pooled-CAV to perform a weighted sum
 188 of the normalized attributions ($A^{t,norm}$). Before this summation, we apply a ReLU and $[0,1]$ rescaling
 189 to the Pooled-CAV (p^c) so that we extract gradually less attribution for feature maps that are less
 190 correlated with the concept. The rationale behind using the ReLU is to discard the attribution of
 191 feature maps that show a negative correlation with the concept. In other words, if a certain feature
 192 map is activated by other non-correlated features, we discard its attribution. Finally, as shown in
 193 Equation (4), we obtain the *Concept Attribution* for a concept (c) and a target class (t) by summing
 194 all values of an element-wise multiplication of the weighted attributions and the concept map (M^c),
 195 which is used as a spatial mask. This enables us to discard the attributions of activations related to the
 196 regions within the input image where the concept is not present or was not recognized.

$$ConceptAttribution_{c,t} = \sum_{i,j} M_{ij}^c \cdot \left(\sum_k ReLU(p_k^{c,norm}) \cdot A_k^{t,norm} \right)_{ij} \quad (4)$$

197 The concept attribution is a per-concept metric of importance, meaning that two concepts can have
 198 significantly different attributions even if they are recognized in the same location of the input
 199 image, resulting in similar concept maps. For instance, considering the “zebra” class, the attribution
 200 of the “striped” concept could be significantly different from the attribution of the “fur” concept.
 201 This distinction is achieved by focusing not on per-pixel attributions but on the attributions of the
 202 activations produced by the neurons responsible for recognizing these two concepts. Moreover, since
 203 the attribution of a concept is independent of its location, we can average it across multiple input
 204 images to provide a quantitative measure of the overall importance of that concept for that particular
 205 class, thus providing a global explanation. For instance, we can calculate a global attribution of the
 206 “striped” concept for the “zebra” target class by averaging the attribution of “striped” across a large
 207 number (e.g., 200) of images containing zebras.

208 4 Experiments and Results

209 In this section, we present the results of applying Visual-TCAV to the following convolutional
 210 networks pre-trained on the ImageNet [6] dataset: GoogLeNet [26], InceptionV3 [27], VGG16 [21],
 211 and ResNet50V2 [9]. Examples of “striped”, “zigzagged”, “waffled”, and “chequered” concepts are

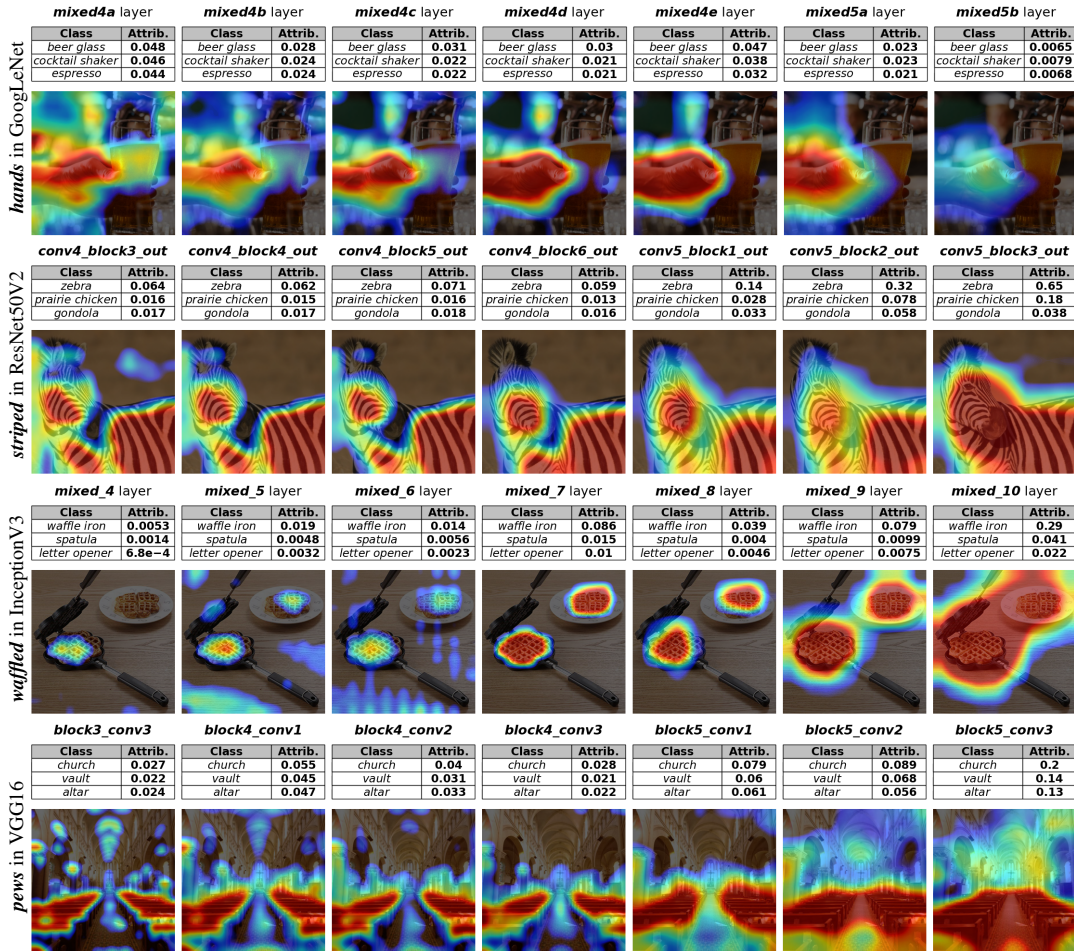


Figure 3: Examples of layer-wise local explanations for various concepts and networks. We compute the attribution of each concept for the top three predicted classes and the last seven layers.

212 sourced from the Describable Textures Dataset (DTD) [4], while “pews” and “fresco” are generated
 213 through Stable Diffusion v1.5 [18] (more on this in Appendix E). Other concepts are obtained from
 214 popular image search engines. Similarly to TCAV, we use a minimum of 30 example images per
 215 concept and 500 random images as negative examples, as suggested by Martin and Weller [15].

216 Our experiments are conducted on an Intel i7 13700k with an Nvidia RTX 4060Ti 16GB, and 32 GB
 217 of DDR5 RAM. The software runs on TensorFlow 2.15.1, CUDA 12.2, and Python 3.11.5. Local
 218 explanations, with 300 steps and seven layers, take less than a minute, while global explanations with
 219 200 class images, 300 steps, and seven layers, can take anywhere from 5 to 20 minutes, depending on
 220 the model. For global explanations, the computation time remains nearly constant regardless of the
 221 number of concepts processed simultaneously. The official implementation is available in our GitHub
 222 repository: *removed for anonymity, see supplemental material .zip file.*

223 4.1 Local Explanations

224 In Figure 3, we provide local explanations for various concepts. While concept maps are class-
 225 independent, the attribution of each concept depends on the class considered. We examine the top
 226 three predicted classes in our examples and apply Visual-TCAV to a subset of the CNNs’ layers. On
 227 one hand, we can observe a substantial increase in attributions in deeper layers, reaching a peak in
 228 the final layer, which holds the most information about the importance of each concept for a specific
 229 class, given its proximity to the output. On the other hand, the most accurate concept maps are
 230 typically found in slightly earlier layers due to their neurons having smaller receptive fields.

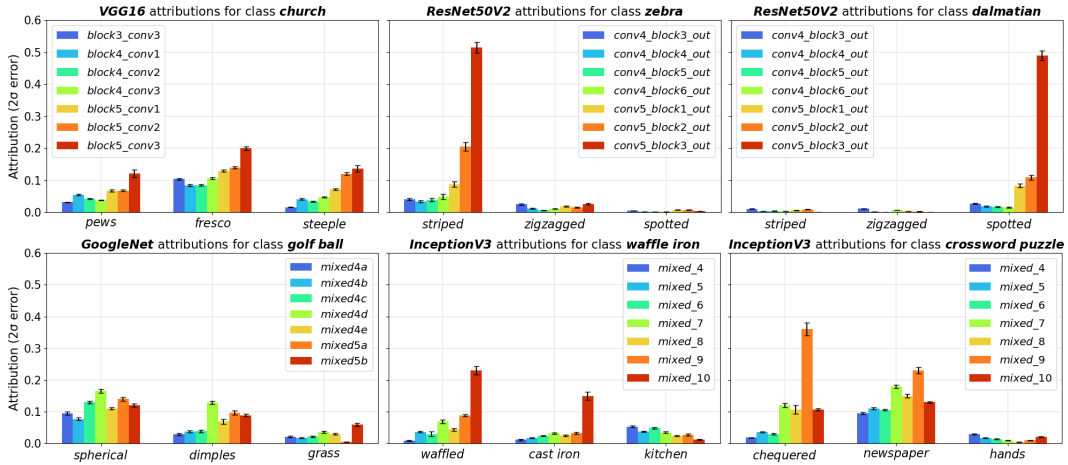


Figure 4: Results of global explanations for a variety of concepts, classes, and networks. Each bar chart reports the attributions of three concepts for a given class, throughout the last seven layers of each network. The attributions of each concept are computed across 200 images of the selected class. Although the theoretical limit of concept attributions is 1.0, the scale in our charts only extends to 0.6. This is based on our empirical observations, which rarely identified concepts with a global attribution exceeding this value.

231 Furthermore, these layer-wise explanations enable us to identify when specific concepts are recognized
 232 within the network. For instance, the “waffled” concept does not significantly activate the initial
 233 layers of InceptionV3, but it is recognized by deeper layers with a considerable attribution in the final
 234 one. We also observe that the “hands” concept is detected mainly by earlier layers and contributes
 235 only marginally to the score of the top classes for the analyzed image. This observation aligns
 236 with the common intuition that “hands” are not class-discriminative in this particular case for the
 237 classes “beer glass”, “cocktail shaker”, and “espresso”. In contrast, the “striped” and “pews” concepts
 238 significantly activate the final layer and substantially contribute to the predictions, although with
 239 different magnitudes of importance. In the case of the “zebra” image, for instance, the network’s
 240 decision is largely influenced by the “striped” concept, which accounts for more than half the logit
 241 value of the “zebra” class. This concept also has a notable impact on the “prairie chicken” class and a
 242 marginal one on the “gondola” class, probably since gondoliers usually wear striped t-shirts. More
 243 examples of local explanations can be found in Appendix C.

244 4.2 Global Explanations

245 The concept attribution is a per-concept metric of importance, hence we can derive global explanations
 246 by aggregating this attribution across a wide range of input images of a selected class. In our
 247 experiments, we utilize 200 images per class for each global explanation. For concepts that are
 248 inherently part of the class (e.g., “striped” for “zebra” or “dimples” for “golf ball”), we can directly
 249 use any image representing that class. On the other hand, for concepts that appear sporadically, we
 250 only use images where the concept is present. For instance, we only use images of church interiors
 251 for “pews” and “fresco” concepts, and images of church exteriors for the “steeple” concept. This
 252 ensures that the explanations are independent of the frequency of the concept’s appearance in the
 253 class images.

254 The results are shown in Figure 4. The attributions match our intuitive expectations, considering, for
 255 instance, the importance of the “striped” concept for “zebra” or “spotted” for “dalmatian”. Moreover,
 256 the final layer typically provides the highest attribution, which is expected for class discriminative
 257 concepts. However, there are instances, such as “chequered” and “newspaper” for “crossword puzzle”,
 258 where concepts recognized in the earlier layers have a greater impact on the network’s prediction. We
 259 observe a more gradual increase in attribution in VGG16 and GoogleNet, compared to InceptionV3
 260 and ResNet50V2. This could be attributed to the depth of the latter networks, which means they
 261 perform more convolution operations that could potentially lead to a more complex feature extraction
 262 between the analyzed layers. More examples of global explanations are provided in Appendix D.

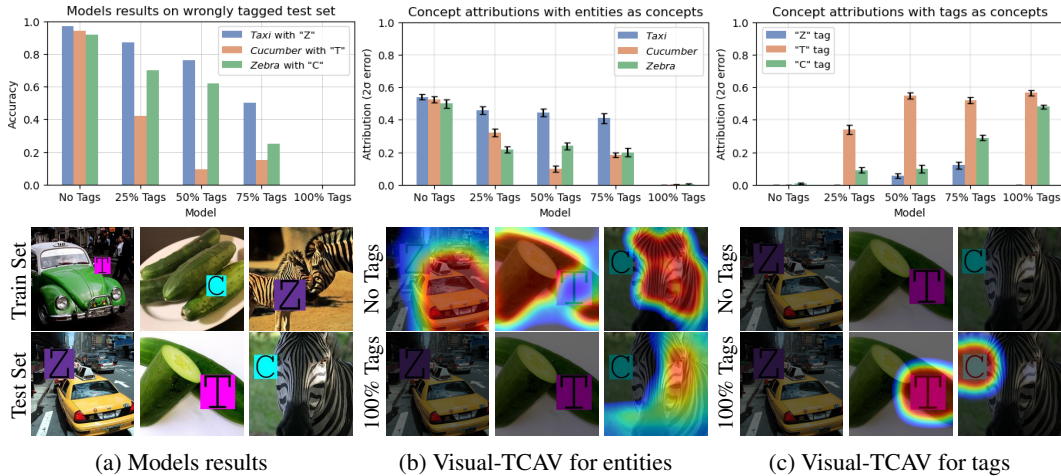


Figure 5: The results of the validation experiment. The upper section of the figure shows the test results and the concept attributions for both entities and tags across all models. The lower section provides examples of tagged images and concept maps for the no tags model and 100% tags model.

263 4.3 Validation Experiment with Ground Truth

264 We conduct a validation experiment to evaluate the effectiveness of Visual-TCAV. In this experiment,
 265 we train convolutional networks in a controlled setting, where ground truth is known, and assess
 266 whether the Visual-TCAV attributions match this ground truth. For this purpose, we create a dataset
 267 of three classes – cucumber, taxi, and zebra – which are the same classes used in the TCAV paper.
 268 We then create multiple versions of this dataset by altering a percentage of the images with a tag,
 269 represented by a letter enclosed in a randomly sized square and added in a random location of the
 270 image (examples are shown in Figure 5a). Specifically, zebra images are tagged with a “Z” in a
 271 purple square, taxi images with a “T” in a magenta square, and cucumber images with a “C” in a
 272 cyan square. From these tagged images, we create five datasets: one of images without tags, and four
 273 others with 25%, 50%, 75%, and 100% of tagged images, respectively. Each dataset is then used
 274 to train a different model, each including six convolutional layers and a GAP layer. Depending on
 275 the dataset used for training, each model may learn to recognize either the entities (i.e., cucumbers,
 276 taxis, and zebras), the tags, or both and will decide which ones to give more importance. To obtain
 277 an approximated ground truth assessing which concept – entity or tag – is more important, we ask
 278 the models to classify a set of 200 incorrectly tagged test images per class. In this test set, taxis are
 279 tagged with the “Z”, cucumbers are tagged with the “T” and zebras are tagged with the “C”. If the
 280 network correctly classifies most of the images, it indicates that the entity is more important than the
 281 tag, and thus, its attribution should be higher. On the other hand, if the performance deteriorates on
 282 these wrongly tagged images, it indicates that the tag is more important than the entity, and thus its
 283 attribution should be higher. We obtain the CAVs for entities using images of each class as concept
 284 examples and random images as negative examples. For tags, we use random images containing that
 285 tag as concept examples and images of cucumbers, taxis, and zebras containing the other two tags as
 286 negative examples. We use the same incorrectly tagged test set to compute the concept attributions
 287 for both entities and tags across the last convolutional layer of all models.

288 The results are shown in Figure 5. As expected, an increase in the percentage of tagged images
 289 correlates with a decrease in accuracy. In particular, for the “cucumber” class the accuracy declines
 290 much faster compared to other classes, with the majority of the images being incorrectly classified
 291 as taxis. This suggests that even the models trained on a small fraction of tagged images tend to
 292 overfit on the “T” tag. The concept attributions for both the “cucumber” entity and the “T” tag
 293 closely mirror this ground truth. The “zebra” entity and the “C” tag are also consistent with the
 294 ground truth: the attributions for “zebra” show a positive correlation with accuracy, whereas the
 295 attributions for the “C” tag demonstrate a clear inverse correlation. Notably, the networks did not
 296 pay much attention to the “Z” tag, focusing instead on the absence of the other two tags to classify
 297 zebras. Indeed, the model trained with 100% of images tagged classifies any image without a “C”
 298 or a “T” tag as “zebra”, regardless of whether the “Z” tag is present or not. This is confirmed by

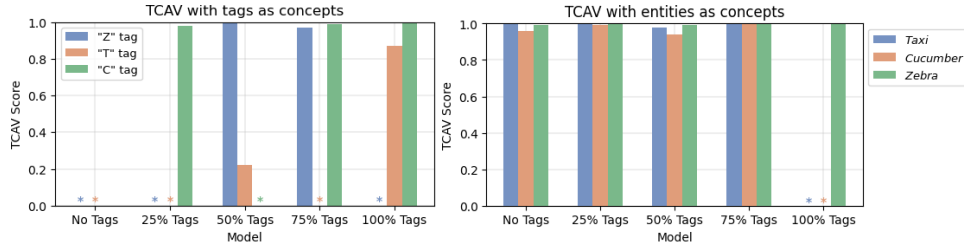


Figure 6: TCAV scores for tags and entities across each validation model. Results marked with an asterisk (“*”) have been excluded due to statistical insignificance (p -value > 0.05).

our method, which assigns an attribution of nearly zero to both the “Z” tag and the “taxi” entity for the aforementioned model. We tested other saliency methods, such as Grad-CAM and IG, to further validate these findings. These methods do not highlight the “Z” tag either, but rather the entire image, in search of the “zebra” class (see Appendix B). For models trained with less than 100% of tags, the accuracy for “taxi” remains high, implying that these models are indeed capable of recognizing the “taxi” entity. The concept attribution for the “taxi” entity aligns with this observation. In Figures 5b and 5c, we provide examples of concept maps for the model trained without tags and the model trained with 100% of tagged images. The former recognizes the entities but not the tags, while the latter struggles to recognize the entities but effectively identifies the “T” and “C” tags.

Comparison with the TCAV Score. The primary difference between our concept attribution and the TCAV score is that the former considers not only the direction of gradients but also their magnitude. This allows us to measure the concept’s impact on the predictions, beyond just the network’s sensitivity to it. To demonstrate this, we compute the TCAV scores for tags and entities across each validation model (see Figure 6). On one hand, TCAV scores match the ground truth in showing that the network trained without tags exhibits high sensitivity to the entities and no sensitivity to the tags. Furthermore, TCAV aligns with the concept attribution in showing that the 100% tags model is sensitive to the “T” and “C” tags but not to the “Z”. On the other hand, TCAV struggles to capture the variations in the concept’s importance defined by ground truth. In fact, all models except the 100% tags show very similar TCAV scores for the entity concepts, even though their importance varies significantly across these models. This is attributable to most of the networks being sensitive to the entities. Indeed, on images without tags, the models’ accuracies are 96.5%, 96.2%, 96.2%, 95.2%, and 36.2% respectively. Similarly, the “C” tag has almost the same TCAV score for the models trained with 25%, 75%, and 100% tags, which is inconsistent with the decline in accuracy for the “C” tagged zebras.

5 Conclusion

In this article, we introduced a novel method, Visual-TCAV, to explain the outputs of image classification models. This framework is capable of providing both local and global explanations based on high-level concepts, by estimating their attribution to the network’s predictions. Additionally, Visual-TCAV generates saliency maps to show where concepts are identified by the network, thereby assuring the user that the attributions correspond to the intended concepts. The effectiveness of this method was demonstrated across a range of widely used CNNs and through a validation experiment, where Visual-TCAV successfully identified the most important concept in each examined model.

Limitations and Future Work. Visual-TCAV provides a novel approach for concept-based explainability, but it has some limitations. Our current implementation only considers positive attributions for classes with positive logit values. However, since a concept may negatively impact the output, in future implementations we aim to include negative values, which would improve explanations and also extend the applicability of Visual-TCAV beyond classification tasks. Another limitation arises from the accumulation of noise along the IG linear path, which may sometimes result in slightly underestimated attributions. Future studies could investigate how to mitigate this using alternative IG variants to compute the attributions of feature maps. Additionally, future research could explore generative approaches such as DreamBooth [19] to generate a large number of concept images starting from a small set of examples, leading to more robust CAVs and reducing workload for analysts. Finally, future works could study interconnections between concepts to determine how the activation of a concept might influence not only the output but also the activation of other concepts.

342 **References**

- 343 [1] Matteo Bianchi, Antonio De Santis, Andrea Tocchetti, and Marco Brambilla. Interpretable
344 network visualizations: A human-in-the-loop approach for post-hoc explainability of cnn-based
345 image classification, 2024. URL <https://doi.org/10.48550/arXiv.2405.03301>.
- 346 [2] Carrie J. Cai, Jonas Jongejan, and Jess Holbrook. The effects of example-based explana-
347 tions in a machine learning interface. In *Proceedings of the 24th International Conference*
348 *on Intelligent User Interfaces*, page 258–26. Association for Computing Machinery, 2019.
349 ISBN 9781450362726. doi: 10.1145/3301275.3302289. URL [https://doi.org/10.1145/](https://doi.org/10.1145/3301275.3302289)
350 [3301275.3302289](https://doi.org/10.1145/3301275.3302289).
- 351 [3] Lei Cai, Jingyang Gao, and Di Zhao. A review of the application of deep learning in medical
352 image classification and segmentation. *Annals of Translational Medicine*, 8(11), 2020. ISSN
353 2305-5847. URL <https://atm.amegroups.com/article/view/36944>.
- 354 [4] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi. Describing textures in the wild.
355 In *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- 356 [5] Jonathan Crabbé and Mihaela van der Schaar. Concept activation regions: A gen-
357 eralized framework for concept-based explanations. In Sanmi Koyejo, S. Mohamed,
358 A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural In-*
359 *formation Processing Systems 35: Annual Conference on Neural Information Process-*
360 *ing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - Decem-*
361 *ber 9, 2022*, 2022. URL [http://papers.nips.cc/paper_files/paper/2022/hash/](http://papers.nips.cc/paper_files/paper/2022/hash/11a7f429d75f9f8c6e9c630aeb6524b5-Abstract-Conference.html)
362 [11a7f429d75f9f8c6e9c630aeb6524b5-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/11a7f429d75f9f8c6e9c630aeb6524b5-Abstract-Conference.html).
- 363 [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-
364 scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern*
365 *Recognition*, pages 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- 366 [7] Amirata Ghorbani, James Wexler, James Y Zou, and Been Kim. Towards automatic concept-
367 based explanations. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and
368 R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran
369 Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2019/file/77d2afcb31f6493e350fca61764efb9a-Paper.pdf)
370 [2019/file/77d2afcb31f6493e350fca61764efb9a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/77d2afcb31f6493e350fca61764efb9a-Paper.pdf).
- 371 [8] Mara Graziani, Vincent Andrearczyk, and Henning Müller. Regression concept vectors for
372 bidirectional explanations in histopathology. In Danail Stoyanov, Zeike Taylor, Seyed Mostafa
373 Kia, Ipek Oguz, Mauricio Reyes, Anne Martel, Lena Maier-Hein, Andre F. Marquand, Edouard
374 Duchesnay, Tommy Löfstedt, Bennett Landman, M. Jorge Cardoso, Carlos A. Silva, Sergio
375 Pereira, and Raphael Meier, editors, *Understanding and Interpreting Machine Learning in*
376 *Medical Image Computing Applications*, pages 124–132, Cham, 2018. Springer International
377 Publishing. ISBN 978-3-030-02628-8.
- 378 [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual
379 networks. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer*
380 *Vision – ECCV 2016*, pages 630–645, Cham, 2016. Springer International Publishing. ISBN
381 [978-3-319-46493-0](https://doi.org/10.1007/978-3-319-46493-0).
- 382 [10] Andrei Kapishnikov, Subhashini Venugopalan, Besim Avci, Ben Wedin, Michael Terry, and
383 Tolga Bolukbasi. Guided integrated gradients: an adaptive path method for removing noise.
384 In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages
385 5048–5056, 2021. doi: 10.1109/CVPR46437.2021.00501.
- 386 [11] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas,
387 and Rory sayres. Interpretability beyond feature attribution: Quantitative testing with concept
388 activation vectors (TCAV). In Jennifer Dy and Andreas Krause, editors, *Proceedings of the*
389 *35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine*
390 *Learning Research*, pages 2668–2677. PMLR, 10–15 Jul 2018. URL [https://proceedings.](https://proceedings.mlr.press/v80/kim18d.html)
391 [mlr.press/v80/kim18d.html](https://proceedings.mlr.press/v80/kim18d.html).

- 392 [12] Zachary Lipton. The mythos of model interpretability. *Communications of the ACM*, 61, 10
393 2016. doi: 10.1145/3233231.
- 394 [13] Adriano Lucieri, Muhammad Naseer Bajwa, Stephan Braun, Muhammad Imran Malik, Andreas
395 Dengel, and Sheraz Ahmed. On interpretability of deep learning based skin lesion classifiers
396 using concept activation vectors. pages 1–10, 07 2020. doi: 10.1109/IJCNN48605.2020.
397 9206946.
- 398 [14] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In
399 I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett,
400 editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates,
401 Inc., 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/file/
402 8a20a8621978632d76c43dfd28b67767-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/8a20a8621978632d76c43dfd28b67767-Paper.pdf).
- 403 [15] Tyler Martin and Adrian Weller. *Interpretable Machine Learning*. M.Phil. diss., Dept. of
404 Engineering, University of Cambridge, August 2019. URL [https://www.mlmi.eng.cam.ac.
405 uk/files/tam_final_reduced.pdf](https://www.mlmi.eng.cam.ac.uk/files/tam_final_reduced.pdf).
- 406 [16] Deng Pan, Xin Li, and Dongxiao Zhu. Explaining deep neural network models with adversarial
407 gradient integration. In Zhi-Hua Zhou, editor, *Proceedings of the Thirtieth International
408 Joint Conference on Artificial Intelligence, IJCAI-21*, pages 2876–2883. International Joint
409 Conferences on Artificial Intelligence Organization, 8 2021. doi: 10.24963/ijcai.2021/396. URL
410 <https://doi.org/10.24963/ijcai.2021/396>. Main Track.
- 411 [17] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. “why should i trust you?”: Explaining
412 the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International
413 Conference on Knowledge Discovery and Data Mining, KDD ’16*. ACM, August 2016. doi:
414 10.1145/2939672.2939778. URL <http://dx.doi.org/10.1145/2939672.2939778>.
- 415 [18] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis
416 with latent diffusion models. In *2022 IEEE/CVF Conference on Computer Vision and Pattern
417 Recognition (CVPR)*, pages 10674–10685, Los Alamitos, CA, USA, jun 2022. IEEE Computer
418 Society. doi: 10.1109/CVPR52688.2022.01042. URL [https://doi.ieeecomputersociety.
419 org/10.1109/CVPR52688.2022.01042](https://doi.ieeecomputersociety.org/10.1109/CVPR52688.2022.01042).
- 420 [19] N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein, and K. Aberman. Dreambooth: Fine tuning
421 text-to-image diffusion models for subject-driven generation. In *2023 IEEE/CVF Conference
422 on Computer Vision and Pattern Recognition (CVPR)*, pages 22500–22510, Los Alamitos,
423 CA, USA, jun 2023. IEEE Computer Society. doi: 10.1109/CVPR52729.2023.02155. URL
424 <https://doi.ieeecomputersociety.org/10.1109/CVPR52729.2023.02155>.
- 425 [20] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi
426 Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based
427 localization. In *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE, October
428 2017. doi: 10.1109/iccv.2017.74. URL <http://dx.doi.org/10.1109/ICCV.2017.74>.
- 429 [21] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale
430 image recognition. In *International Conference on Learning Representations*, 2015.
- 431 [22] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks:
432 Visualising image classification models and saliency maps. In *Workshop at International
433 Conference on Learning Representations*, 2014.
- 434 [23] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda B. Viégas, and Martin Wattenberg.
435 Smoothgrad: removing noise by adding noise. *CoRR*, abs/1706.03825, 2017. URL [http:
436 //arxiv.org/abs/1706.03825](http://arxiv.org/abs/1706.03825).
- 437 [24] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin A. Riedmiller. Striving
438 for simplicity: The all convolutional net. *CoRR*, abs/1412.6806, 2014.
- 439 [25] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In
440 *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML’17*,
441 page 3319–3328. JMLR.org, 2017.

- 442 [26] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov,
 443 Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions.
 444 In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–9,
 445 2015. doi: 10.1109/CVPR.2015.7298594.
- 446 [27] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Re-
 447 thinking the inception architecture for computer vision. In *2016 IEEE Conference on Computer
 448 Vision and Pattern Recognition (CVPR)*, pages 2818–2826, 2016. doi: 10.1109/CVPR.2016.308.
- 449 [28] Tolga Turay and Tanya Vladimirova. Toward performing image classification and object
 450 detection with convolutional neural networks in autonomous driving systems: A survey. *IEEE
 451 Access*, 10:14076–14119, 2022. doi: 10.1109/ACCESS.2022.3147495.
- 452 [29] Warren von Eschenbach. Transparency and the black box problem: Why we do not trust ai.
 453 *Philosophy & Technology*, 34, 12 2021. doi: 10.1007/s13347-021-00477-0.
- 454 [30] Chase Walker, Sumit Jha, Kenny Chen, and Rickard Ewetz. Integrated decision gradients:
 455 Compute your attributions where the model makes its decision. *Proceedings of the AAAI
 456 Conference on Artificial Intelligence*, 38:5289–5297, 03 2024. doi: 10.1609/aaai.v38i6.28336.
- 457 [31] Ruihan Zhang, Prashan Madumal, Tim Miller, Krista A. Ehinger, and Benjamin I. P. Rubinstein.
 458 Invertible concept-based explanations for cnn models with non-negative concept activation
 459 vectors. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(13):11682–11690,
 460 May 2021. doi: 10.1609/aaai.v35i13.17389. URL [https://ojs.aaai.org/index.php/
 461 AAAI/article/view/17389](https://ojs.aaai.org/index.php/AAAI/article/view/17389).

462 A Appendix Overview

463 In the appendix, we provide:

- 464 B. Saliency methods for 100% tags model
- 465 C. Additional results of Local Explanations
- 466 D. Additional results of Global Explanations
- 467 E. Example images for generated concepts

468 B Saliency methods for 100% tags model

469 We provide the results obtained by applying IG and Grad-CAM to the 100% tags model (see Figure 7).
 470 These methods align with Visual-TCAV in showing that this model does not pay attention to the “Z”,
 471 but rather to the absence of the “T” and the “C” for predicting the “zebra” class.

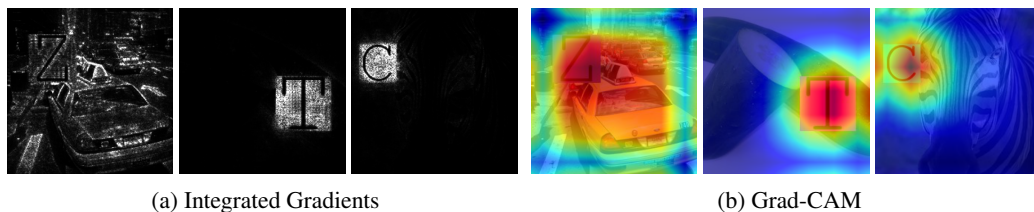


Figure 7: Integrated Gradients and Grad-CAM for the model with 100% tags, searching respectively for the classes “zebra”, “taxi”, and “cucumber”. Both methods highlight the “T” for class “taxi” and the “C” for class “cucumber”, but fail to recognize the “Z” for class “zebra”.

472 C Additional results of Local Explanations

473 Continuing from the results presented in Section 4.1, we further provide additional local explanations
 474 for more input images and concepts in Figure 8.

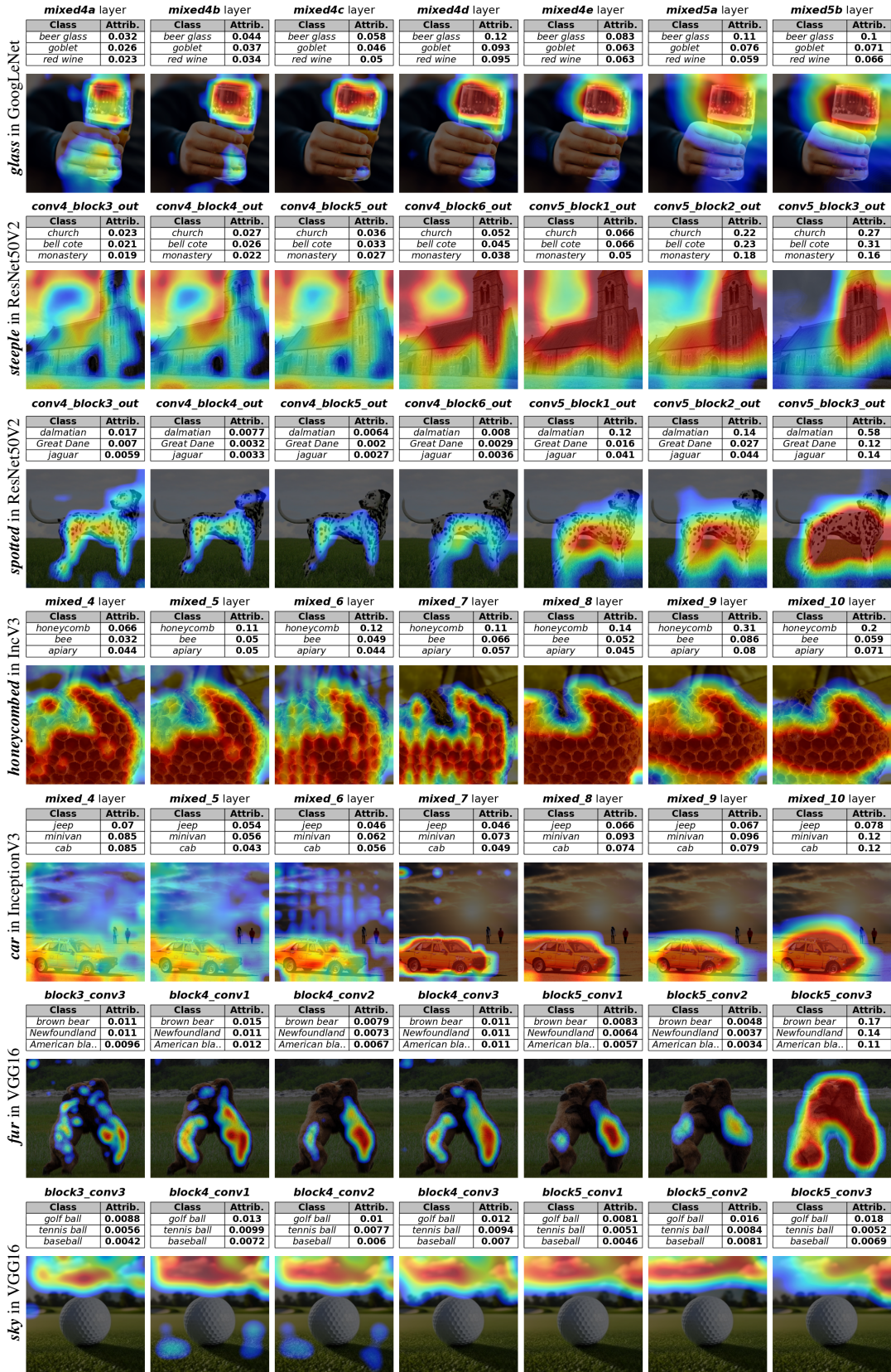


Figure 8: More examples of layer-wise local explanations for various concepts and networks.

475 **D Additional results of Global Explanations**

476 Building upon the results outlined in Section 4.2, we provide additional global explanations for
 477 various classes and concepts in Figure 9.

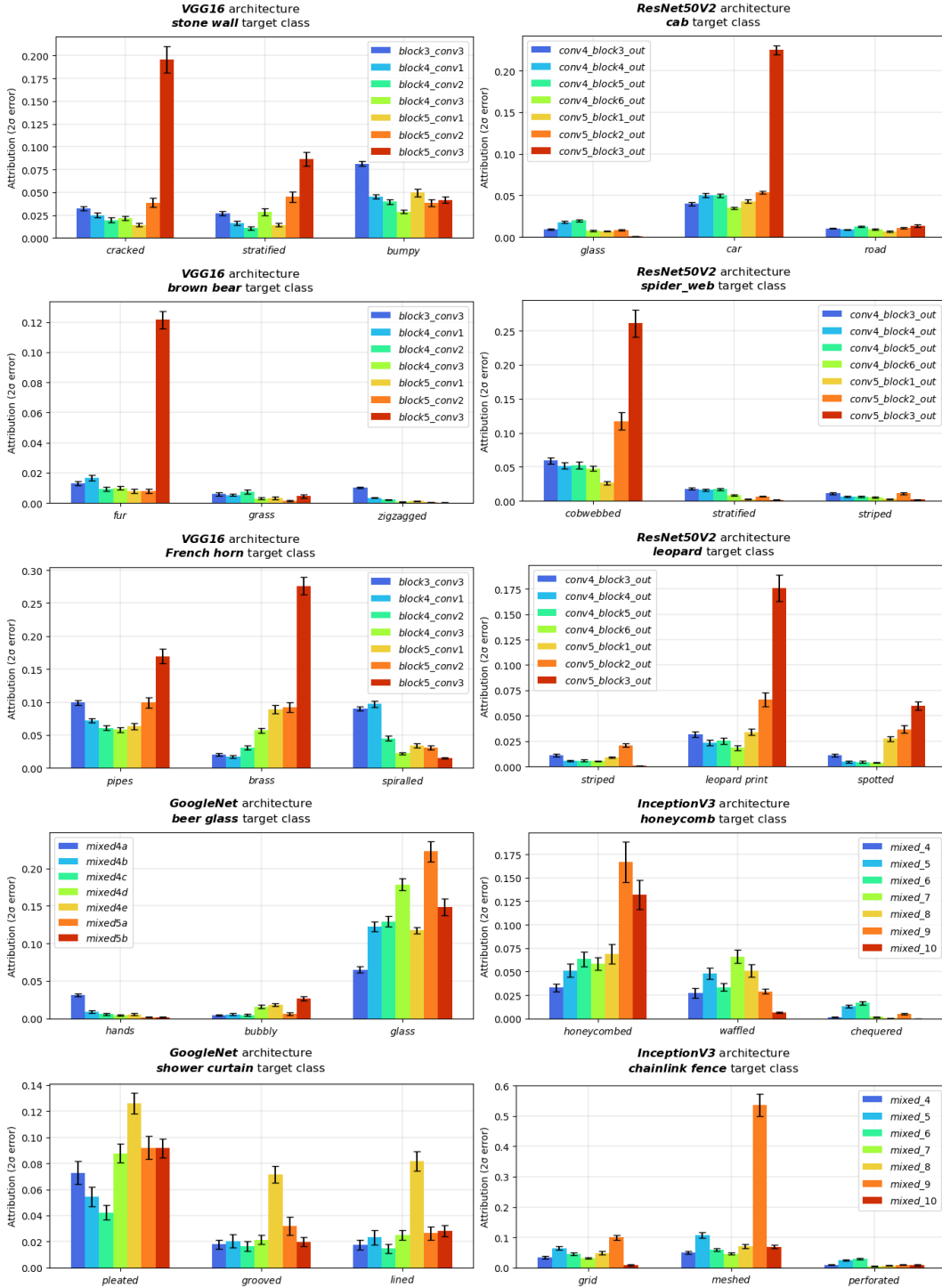


Figure 9: More examples of global explanations for various classes, concepts, and networks.

478 **E Example images for generated concepts**

479 Some of the concepts used in the paper were automatically generated using Stable Diffusion v1.5 [18]
480 with default parameters. In particular, we generated the following concepts: “pews”, “fresco”,
481 “arches”, “sky”, “pipes”, and “brass”. We used just the concept name as a prompt and generated 200
482 images per concept. A subsequent manual revision was still necessary to eliminate errors and strange
483 artifacts. In Figure 10, we provide three example images for each generated concept.

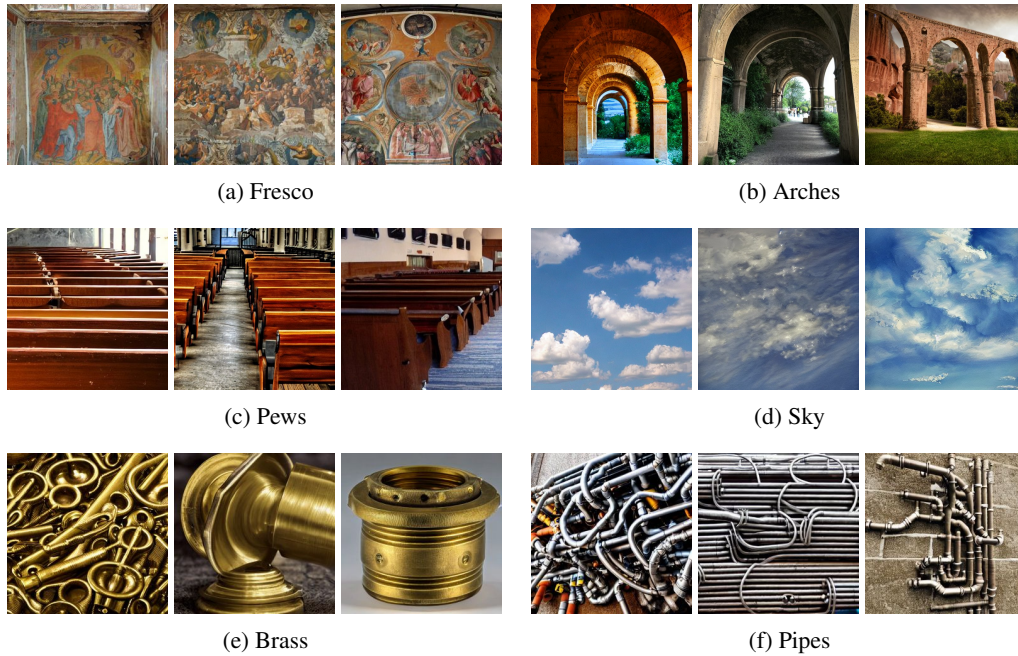


Figure 10: We provide three example images for each concept generated with Stable Diffusion v1.5.

484 **NeurIPS Paper Checklist**

485 **1. Claims**

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

488 Answer: [Yes]

489 Justification: We claim that our method can provide visual explanations through saliency
490 maps based on user-defined concepts, estimate the attributions of these concepts for a
491 selected class, and provide both local and global explanations. These claims are all validated
492 through the experimental results performed in the paper.

493 **2. Limitations**

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

495 Answer: [Yes]

496 Justification: Our work has some limitations, we discuss them in Section 5.1.

497 **3. Theory Assumptions and Proofs**

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

500 Answer: [NA]

501 Justification: The paper does not include any new proof or theorem.

502 **4. Experimental Result Reproducibility**

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

506 Answer: [Yes]

507 Justification: Every experimental result presented in the paper is fully reproducible using
508 the provided code and data.

509 **5. Open access to data and code**

510 Question: Does the paper provide open access to the data and code, with sufficient instruc-
511 tions to faithfully reproduce the main experimental results, as described in supplemental
512 material?

513 Answer: [Yes]

514 Justification: We provide the code, data, and instructions needed to reproduce every ex-
515 periment both to reviewers and to the public through a GitHub repository (in case of
516 publication).

517 **6. Experimental Setting/Details**

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

521 Answer: [Yes]

522 Justification: The paper describes in detail all the necessary steps to reproduce and un-
523 derstand the experiments. Furthermore, the code used is also available as supplementary
524 material.

525 **7. Experiment Statistical Significance**

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

528 Answer: [Yes]

529 Justification: In our bar plots we always report 2-sigma error bars.

530 **8. Experiments Compute Resources**

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

534 Answer: [Yes]

535 Justification: We describe in detail the characteristics of the machine used to run all the
536 experiments and the execution time.

537 9. Code Of Ethics

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

540 Answer: [Yes]

541 Justification: We have reviewed the NeurIPS Code of Ethics and our research conforms with
542 it.

543 10. Broader Impacts

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

546 Answer: [Yes]

547 Justification: In the introduction, we briefly discuss the problem of transparency in AI
548 systems, particularly as Convolutional Neural Networks are being widely utilized in critical
549 sectors such as healthcare and autonomous driving. Our work can have a positive societal
550 impact by facilitating a trustworthy adoption of these systems. We are not aware of any
551 negative impact our work could have.

552 11. Safeguards

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

556 Answer: [NA]

557 Justification: This paper does not release any data or models that pose such risks.

558 12. Licenses for existing assets

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

562 Answer: [Yes]

563 Justification: All models and datasets used for the experiments are properly cited in the
564 paper.

565 13. New Assets

566 Question: Are new assets introduced in the paper well documented and is the documentation
567 provided alongside the assets?

568 Answer: [NA]

569 Justification: The paper does not introduce new assets.

570 14. Crowdsourcing and Research with Human Subjects

571 Question: For crowdsourcing experiments and research with human subjects, does the paper
572 include the full text of instructions given to participants and screenshots, if applicable, as
573 well as details about compensation (if any)?

574 Answer: [NA]

575 Justification: The paper does not involve crowdsourcing nor research with human subjects.

576 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 577 Subjects

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

582 Answer: [NA]

583 Justification: The paper does not involve crowdsourcing nor research with human subjects.