

FROM ATTENTION TO PREDICTION MAPS: PER-CLASS GRADIENT-FREE TRANSFORMER EXPLANATIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

The Vision Transformer (ViT) has become a standard model architecture in computer vision, especially for classification tasks. As such, explaining ViT predictions has attracted significant research efforts in recent years. Many methods rely on attention maps, which highlight *where* in the image the network directs its attention. In this paper, we introduce Prediction Maps – a novel explanation method that complements attention maps by revealing *what* the network sees. Prediction maps visualize how each patch token within a given layer is associated with each possible class. This is done by utilizing the classification head at the output of the network, originally trained to be fed with the class token at the last layer. Specifically, to obtain the prediction map of a particular layer, we apply the classification head to every patch token within that layer. We show that prediction maps provide complementary information to attention maps and illustrate that combining them leads to state-of-the-art explainability performance. Furthermore, since our proposed method is neither gradient- nor perturbation-based, it offers superior computational and memory efficiency compared to competing methods. To the best of our knowledge, ours is the first explainability method for ViTs that is both class-specific and gradient-free

1 INTRODUCTION

Following their introduction in the context of language models, transformers (Vaswani et al., 2017) have become the neural architecture of choice across diverse machine learning domains. They have been adopted e.g. in graph neural networks (Dwivedi & Bresson, 2020; Yun et al., 2019) and for point-cloud analysis (Qin et al., 2022; Zhao et al., 2021), and have also been extended to a wide range of vision tasks, including detection (Carion et al., 2020; Li et al., 2022; Misra et al., 2021), classification (Dosovitskiy et al., 2020; Zhao et al., 2021), segmentation (Kirillov et al., 2023; Zheng et al., 2021), and image generation (Li et al., 2019; Touvron et al., 2023).

Given their pervasive dominance, significant research efforts have been devoted to understanding how transformers process their inputs, as well as to explaining their predictions, with particular focus given to the vision transformer (ViT) architecture (Abnar & Zuidema (2020); Chefer et al. (2021a;b); Liu et al. (2021a); Mohankumar et al. (2020); Wu et al. (2024)). Many methods rely on the attention maps within the model to provide explanations for its predictions (Abnar & Zuidema, 2020; Chefer et al., 2021b). Raw attention maps are appealing because (i) they are calculated as part of the forward-pass of the network and thus do not require any additional computations to extract, and (ii) they provide a glimpse into how the model constructs its prediction. However, attention maps only offer insights into *where* in the input image the network focuses its attention, and do not visualize *what* the network “perceives” within each region of the image. In other words, they do not indicate the extent to which each region is associated with a specific class.

A common approach for providing more informative explanations, is to seek for heatmaps that visualize the contribution of each patch in the input image to each possible class prediction. Methods that generate such visualizations can be broadly categorized as *perturbation based* (Carter et al., 2019; Fong et al., 2019; Fong & Vedaldi, 2017; Lundberg & Lee, 2017; Petsiuk et al., 2018; Ribeiro et al., 2016) or *gradient based* (Bach et al., 2015; Chefer et al., 2021a;b; Selvaraju et al., 2017; Sundararajan et al., 2017). Perturbation-based methods treat the model as a black box, inspecting how its output changes in response to small perturbations to its input. Gradient-based methods

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

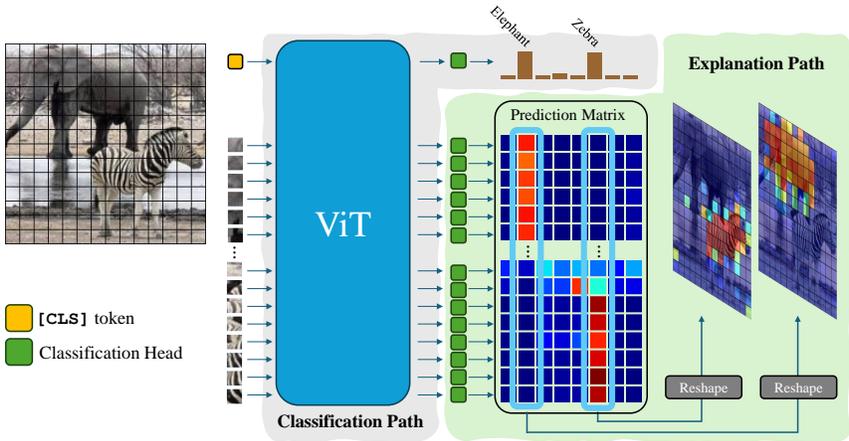


Figure 1: **Prediction map construction overview.** Applying the classification head on all patch tokens yields a class-specific per token classification. Although it was originally trained on the class token, the predictions obtained from other patches are satisfying.

perform a backward pass to accumulate gradients propagated through the entire network. However, both approaches are associated with heavy computational costs, and do not directly shed light on how the model processes data in its forward pass.

In this paper, we introduce *Prediction Maps* – a lightweight gradient-free explainability method that is as simple and fast as extracting attention maps and has the per-class expressiveness of the sophisticated perturbation- and gradient-based methods. Our approach relies on the observation that when the classification head of a pretrained ViT, which is normally applied to the class-token at the last layer, is fed with any other (patch) token, it tends to output valid predictions. Surprisingly, we find that this is true not only for the patch tokens at the last layer, but also for the patch tokens at all other layers. This observation allows us to construct a heatmap for any desired class, as demonstrated in Fig. 1, as well as to visualize how the localization of concepts evolves throughout the layers.

Prediction maps provide complementary information to attention maps. Therefore, their joint inspection can yield a more comprehensive explanation than each of them alone, enabling to expose the root cause of incorrect predictions. To illustrate this, we show in Fig. 2 two failure cases of a ViT on images from the ImageNet-R dataset (Hendrycks et al., 2021a). In the first, the model misclassifies a lemon as tray. Here, the attention and prediction maps reveal that the lemon is actually recognized correctly, but is not attended by the model. In the second example, the model misclassifies a violin as a hair slide. Here, the attention and prediction maps reveal that the violin is well attended, yet not recognized, possibly because of the pencil strokes that look like hair.

To obtain a single map that combines the explainability power of both approaches, we propose to unify them into a visualization which we term *PredicAtt*. We use the fact that the correlation between an attention map and a prediction map indicates how much that attention map contributes to the classification. We therefore construct a weighted combination of the attention maps from all heads and layers according to their similarity with the prediction map. We then compute the per-element product between this weighted attention map and the prediction map to obtain our combined class-specific map. This enables the analysis of how each layer and head recognizes each class.

Our main contributions can be summarized as follows:

1. We introduce *prediction maps* for explainability as a complementary component for the well studied attention maps. Prediction maps provide a per-token measurement of what the network perceives. Additionally, we propose a simple way to combine attention maps and prediction maps, termed *PredicAtt*.
2. Our method is gradient-free and perturbation-free; hence, it is efficient in terms of runtime and memory consumption. To the best of our knowledge, it is the first explainable method for ViT that is both class-specific and gradient-free.

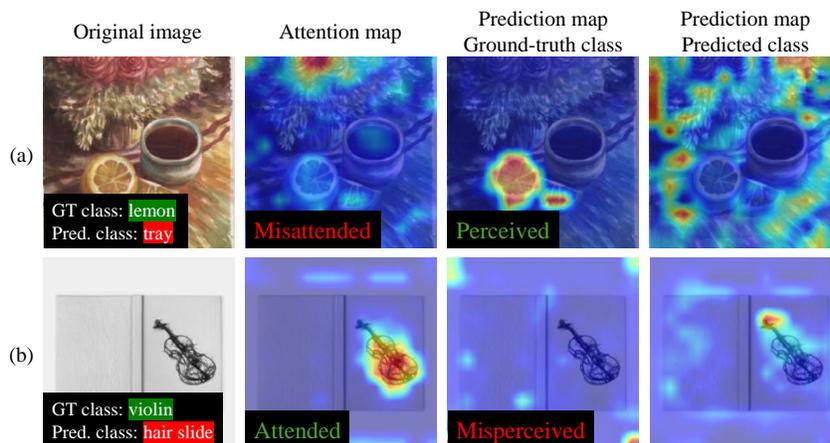


Figure 2: **Gradient-free analysis of failure cases.** (a) The lemon object is perceived accurately, but ViT attends elsewhere. (b) The object is accurately attended; nevertheless, the network mispredicts it, possibly due to the pencil strokes that look like hair. Analyzing ViT through the lens of prediction maps, which complement attention maps, provides interesting insights into the root cause of failure cases. Prediction map of the predicted class is given for completeness.

3. We show that the correlation between prediction maps and attention maps offers insights into how the ViT processes data at a head granularity level. Combining prediction maps with attention maps provides per-class explanations, enhancing interpretability.
4. Our approach achieves state-of-the-art results on perturbation and segmentation tests utilizing ViT-B and ViT-L on the ImageNet dataset (Russakovsky et al., 2015).

2 RELATED WORK

Explainable AI has been widely studied across various domains utilizing deep neural networks, including NLP (Li et al., 2015; Chefer et al., 2021a), speech (Bharadhwaj, 2018; Kumar et al., 2021), point cloud analysis (Levi & Gilboa, 2024; Zheng et al., 2019), and graph neural networks (Ying et al., 2019; Yuan et al., 2021). In the realm of image classifiers, many works aim at producing a heatmap that highlights which regions in the input image affect the classifier’s prediction the most.

Explainability for arbitrary architectures. Generic methods for explaining the predictions of image classification models can be roughly categorized as gradient-based or perturbation-based. Gradient based methods use backpropagation in various ways. For example, GradCAM (Selvaraju et al., 2017) computes the gradient of the score for any queried class with respect to the last pooling layer. Layer-wise Relevance Propagation (LRP) (Bach et al., 2015) propagates relevance values, which are based on gradients, through the entire network. Integrated Gradients (Sundararajan et al., 2017) accumulates the gradients on a path from a baseline image to the scrutinized one. These approaches, however, are computationally demanding. Moreover, they do not shed light on how the computations within the network’s forward pass lead to its prediction. Perturbation-based methods are model-agnostic, treating the model as a black box and inspecting how its output changes in response to small perturbations to its input. LIME (Ribeiro et al., 2016) learns a linear model based on perturbations of the input sample. SHAP (Lundberg & Lee, 2017) translates the concept of Shapley values to model explainability. RISE (Petsiuk et al., 2018) estimates pixel importance based on probing the model with randomly masked versions of the input image. Fong & Vedaldi (2017) optimize perturbations on the input data to identify the smallest most influential regions that significantly affect the model’s output. Extreme Image Transformations (Malik et al., 2023) is a more computationally efficient approach, analyzing how significant alterations to the input influence the model’s predictions. Carter et al. (2019) identify minimal subsets of features whose observed values alone suffice for the same decision to be reached. Despite being applicable to any model, many of these methods suffer from extreme computational costs.

Explainability for convolutional networks. Class Activation Mapping (CAM) (Zhou et al., 2016) generates explainable maps for convolutional neural networks (CNNs) by leveraging the weighted sum of the feature maps from the final convolutional layer, guided by class-specific weights from the fully connected layer. Our prediction maps can be viewed as an adaptation and extension of CAM to ViTs. While CAM in CNNs can only be applied to the last layer due to differing channel dimensions across earlier layers, prediction maps can be applied to any layer in ViTs. Furthermore, we demonstrate how to combine prediction maps with attention maps to enhance the explainability.

Explainability for transformers. With the widespread adoption of transformers (Vaswani et al., 2017), many explainability methods emerged particularly for those architectures. Most methods make use of the self-attention mechanism, utilizing attention maps for explainability (Abnar & Zuidema, 2020; Chefer et al., 2021a;b; Liu et al., 2021a; Mohankumar et al., 2020; Wu et al., 2024). While early works used raw attention maps for explanation, subsequent approaches explored more sophisticated techniques, such as rolling out attention information from all layers (Abnar & Zuidema, 2020). Voita et al. (2019) applied LRP to transformers, focusing only on attention head relevance. Chefer et al. (2021b) further adapted LRP, allowing propagating relevance scores through all layers, while Chefer et al. (2021a) generalized the method to models with cross-attention modules. Wu et al. (2024) incorporated the influence of token transformations into their explainable method, particularly focusing on changes in the tokens’ norms and directions in each attention block.

3 BACKGROUND

In this section we provide a brief overview of the Vision Transformer (ViT) architecture, mainly to set notations. A more comprehensive description can be found in Dosovitskiy et al. (2020).

A ViT consists of a stack of L transformer encoder layers, each comprising a multi-head-self-attention (MHSA) module and a feed-forward network (FFN) block with skip connections (He et al., 2016). The input to each layer is a sequence of tokens, each corresponding to a distinct patch within the input image, along with an additional special token for classification called the [CLS] token (Devlin et al., 2018). The [CLS] token is designed to aggregate information from all other tokens to enable classification based on that token alone. Therefore, at the last layer, the [CLS] token is fed into a classification head, which outputs the predicted class.

We denote by $x^{(l)} \in \mathbb{R}^{(N+1) \times d}$ the output of layer l , where N is the number of regular (patch) tokens (excluding the [CLS] token) and d is the embedding size. Thus, $x^{(0)}$ is the network’s input embedding and $x^{(L)}$ is the output of the last layer. At the output of layer l , we denote the [CLS] token by $x_{\text{CLS}}^{(l)} \in \mathbb{R}^d$ and the i th patch token by $x_i^{(l)} \in \mathbb{R}^d$, for any $i \in \{1, \dots, N\}$. Each transformer layer contains a MHSA block with H heads, each outputting a vector of dimension d_H per token. These H vectors are concatenated and linearly transformed back to dimension d . The h th head within the l th layer applies three linear transformations to $x^{(l)}$, with matrices $\{Q^{(l,h)}, K^{(l,h)}, V^{(l,h)}\} \in \mathbb{R}^{(N+1) \times d_H}$. The results of these transformations are the queries, keys and values, respectively. An *attention matrix* is then constructed as

$$A^{(l,h)} = \text{softmax} \left(\frac{Q^{(l,h)} K^{(l,h)\top}}{\sqrt{d_H}} \right) \in \mathbb{R}^{(N+1) \times (N+1)}, \quad (1)$$

where the softmax operates on the rows of its matrix argument. The *attention map* $A_{\text{CLS}}^{(l,h)} \in \mathbb{R}^N$ is the row of $A^{(l,h)}$ corresponding to the [CLS] token, excluding the entry of the attention between the [CLS] token and itself. The attention map captures the relation between the [CLS] token and each patch token. Reshaping it back to the image dimensions may serve to explain the model’s prediction.

4 LIMITATIONS AND CHALLENGES IN USING ATTENTION MAPS FOR EXPLAINABILITY

The question of whether raw attention maps provide informative explanations has attracted a lot of debate (Bibal et al., 2022; Jain & Wallace, 2019; Wiegrefe & Pinter, 2019). Here, we list several of the obvious limitations and challenges associated with using them for explainability.

Perhaps the most fundamental limitation of attention maps is their class-agnostic nature. Since there is no attention map per class, they do not reveal to what extent each token “recognizes” each possible class. Such a visualization is often highly desirable for explaining a model’s prediction. For example, in cases of incorrect classification, it is informative to inspect the extent to which each token “recognizes” the true class vs. the predicted class (see Fig. 2). Additionally, when the input image contains multiple objects, it is desirable to visualize how each token “recognizes” each of the classes present in the image.

Recent work, such as the Multi-class Token Transformer (MCTformer) (Xu et al., 2022), demonstrates an alternative approach to obtaining class-specific attention maps by introducing multiple class tokens, each representing a distinct class. This approach is both gradient-free and class-specific, addressing the limitations of conventional attention maps. However, unlike MCTformer, our method achieves class-specificity without requiring architectural changes or additional class tokens, offering a more generalizable solution for visualizing class-discriminative behavior in standard Vision Transformers.

Another challenge in using attention maps for explainability relates to the difficulty in determining how to select or combine attention maps from the different layers and heads. Some works use the attention maps of the last layer, assuming they capture the most high-level semantics (Caron et al., 2021). Other methods attempt to find a better combination of attention maps (Abnar & Zuidema, 2020; Chefer et al., 2021a;b). For example, the Rollout (Abnar & Zuidema, 2020) method uses $(I + A^{(1)}) \cdot (I + A^{(2)}) \dots (I + A^{(L)})$, where I is the identity matrix. Yet, determining the optimal combination of attentions from different layers is still an active area of research. As for selecting the attention maps from the different heads within each layer, some methods heuristically opt to average them (Abnar & Zuidema, 2020; Tang et al., 2018; Voita et al., 2018). However, it has been demonstrated that different heads may capture different semantics (Voita et al., 2019), even within the same layer. Recently, Darcet et al. (2023) showed that the network sometimes encodes global information in some of the tokens. In such cases, the corresponding entries in the attention map may no longer represent the original semantic meanings of the tokens, and the effect of averaging across attention heads remains unclear.

Finally, relying solely on attention maps may fail in the face of dataset biases, where the network can exploit shortcuts or spurious cues (Geirhos et al., 2020; Hendrycks et al., 2021b), resulting in the highlighting of non-discriminatory regions. For instance, it has been shown that the attention maps of a cow in a pasture image tend to over-attend to the grass rather than the cow itself. This is despite the fact that the grass is not the most important region for constructing the network’s prediction, as revealed by perturbation tests (Chefer et al., 2022). That is, the network persists to classify the cow correctly even when the grass patches are masked.

5 PREDICTION MAPS

Attention maps capture the interactions between the queries and keys, indicating *where the network looks*. Yet, they do not provide insights into *what the network sees*, which is encoded in the values within the self-attention mechanism. To address this gap, we introduce the concept of *Prediction Maps* which decode the data in the values into a human-understandable visualization. Recall that the classification head, $\nu : \mathbb{R}^d \rightarrow \mathbb{R}^C$, is normally fed with $x_{\text{CLS}}^{(L)}$, where C is the number of classes, to output the prediction $y = \nu(x_{\text{CLS}}^{(L)})$. Here, we propose to feed this head with the patch tokens, rather than the [CLS] token. Specifically, we define the *prediction matrix* of layer l as

$$\Psi^{(l)} := \left[\nu \left(x_1^{(l)} \right) \quad \nu \left(x_2^{(l)} \right) \quad \dots \quad \nu \left(x_N^{(l)} \right) \right]^T \in \mathbb{R}^{N \times C}. \quad (2)$$

This matrix provides a classification result based on each token within layer l separately (see Fig. 1). The i th row of this matrix is the probability vector associated with predicting all possible classes based on the i th token. The c th column, denoted as $\Psi^{(l)}(c) \in \mathbb{R}^N$, forms the *prediction map* of class c . This map contains the prediction scores of class c obtained for each patch token.

Surprisingly, despite the fact that the classification head has been originally trained to operate on the [CLS] token, we illustrate empirically that feeding it with patch tokens yields sensible results (see Fig. 3). This can be attributed to the fact that the [CLS] token and the patch tokens undergo the same processing within the model. Specifically, every token at the output of an attention layer

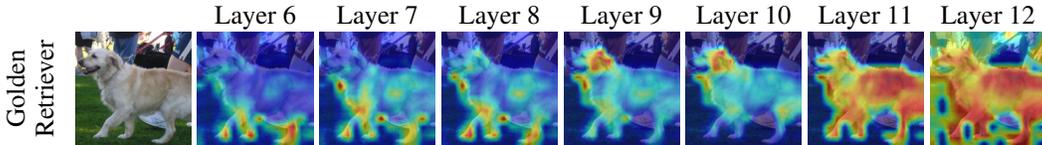


Figure 3: **Prediction map per layer.** Surprisingly, the classification head of a pretrained ViT works well on all tokens from all layers, albeit being originally trained to predict the class only based on the class-token of the last layer. The deeper the layer, the more it captures high-level semantics. We hypothesize that the last layer’s prediction is worse than that of the penultimate layer, because the patch tokens at its output did not participate in the training of the model.

(including the [CLS] token) is a linear combination of (the Values of) all the tokens at the input of that layer. Furthermore, the weights of this linear combination are computed using the Queries and Keys of the tokens at the input, which are computed in the exact same manner for all tokens (namely, the same set of matrices is applied to each token to compute its Q , K , and V). Therefore, there is perfect symmetry in how the tokens are treated, which causes all tokens (including the [CLS] token) to encode information in the same manner. This is why applying the classification head to patch tokens provides sensible classification logits, just like when applying it to the [CLS] token. We note that the object detection method of Minderer et al. (2022) proposed to train a classifier based on patch tokens of the last ViT layer. This is in sharp contrast to our method, which uses the pretrained classification head ‘as is’ and also applies it to tokens from earlier layers.

The prediction map mechanism offers several advantages for explainability. First, it allows to visualize *what the network sees* in each region of the image. Second, it is computationally efficient, as it uses the tokens that are anyway computed during the forward pass of the network. Third, it sheds light on how the network constructs its prediction, as it does not rely on indirect measures based on gradients. Lastly, it provides class-specific maps. As opposed to attention maps, prediction map can be constructed for any desired class. For example, for an ImageNet classifier, it is possible to construct 1,000 prediction maps for each head within each layer, one prediction map per class.

5.1 PREDICATT: COMBINING PREDICTION MAPS WITH ATTENTION MAPS

While prediction maps offer advantages over attention maps, they lack information regarding where the network attends. To incorporate such information, we propose to integrate prediction maps and attention maps into a unified visualization, which we term *PredicAtt*. As we show in Sec. 6, this leads to state-of-the-art explainability results.

To generate a combined map for a specific class c , we follow a two-step process. First, we construct a combined attention map for that class, $\tilde{A}_{\text{CLS}}(c) \in \mathbb{R}^N$, by computing the weighted sum of attention maps across all heads and layers as

$$\tilde{A}_{\text{CLS}}(c) := \sum_{l=1}^L \sum_{h=1}^H \tilde{\alpha}_{l,h} A_{\text{CLS}}^{(l,h)}. \quad (3)$$

The coefficients in this weighted combination are computed based on the similarity between the attention maps and the prediction map of layer i ,

$$\begin{aligned} \alpha_{l,h} &= \left\langle A_{\text{CLS}}^{(l,h)}, \Psi^{(i)}(c) \right\rangle, \\ \{\tilde{\alpha}_{l,h}\} &= \text{softmax}(\{\alpha_{l,h}\}), \end{aligned} \quad (4)$$

where $\langle \cdot, \cdot \rangle$ is the standard inner-product $\langle a, b \rangle = a^\top b$. For details regarding an alternative similarity measure, please refer to Sec. A.2 of the supplementary. In the second step, we compute the element-wise product between the class-specific weighted attention map and the class-specific prediction map, to yield

$$\text{PredicAtt}_i(c) := \tilde{A}_{\text{CLS}}(c) \odot \Psi^{(i)}(c) \in \mathbb{R}^N. \quad (5)$$

In principle, the prediction map $\Psi^{(i)}(c)$ in Eq. (5) can be taken from any layer i . However, as we show in Sec. 6.1, the last and second-to-last layers lead to the best results.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

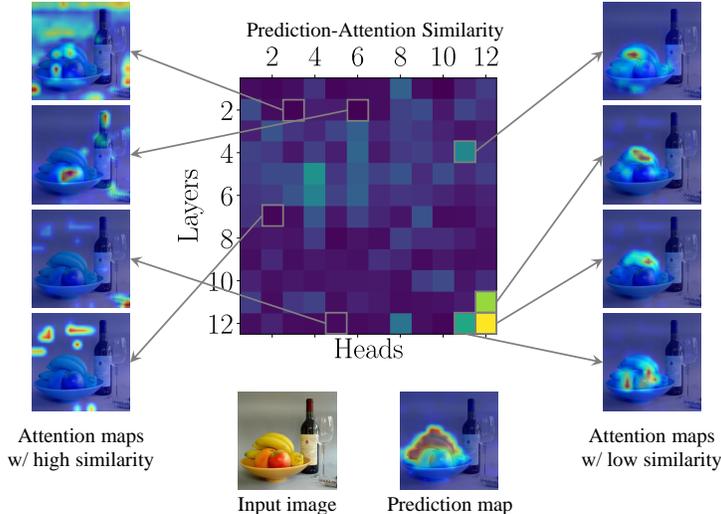


Figure 4: **Analyzing the contribution of attention maps through similarity to prediction maps.** The similarity between an attention map and the prediction map reflects the degree to which this attention map contributes to the prediction. Attention maps with high similarity to the prediction map for class “banana”, emphasize the banana object. In contrast, attention maps with low similarity attend to the background, and are thus less informative for the classification. Interestingly, the least and most contributing maps are not necessarily from the first and last layers, respectively. See the supplementary for an example (Fig. 13).

Our approach is based on a per-head, rather than per-layer analysis and thus facilitates a more detailed understanding of the network’s operation. We claim that an attention map significant for classification is one that yields a high dot product score with the prediction map. A higher dot product score indicates that the network is attending to a relevant region of the image that is mostly associated with a single class. On the other hand, if an attention map resonates over regions containing multiple classes, then its impact on the classification is intuitively smaller, and accordingly, it results in a lower dot product score. In other words, a high degree of similarity between the *what the network sees* template and the *where to look* template signifies that the network has identified meaningful features (see Fig. 4). We empirically verify this in Sec. 6.1 (Tab. 4, rows (i) and (ii)).

6 EXPERIMENTS

We evaluate our method on the task of explaining ViT classifier predictions. In Sections A.4 and A.5 of the supplementary, we further illustrate our approach on explaining the text model BERT (Devlin et al., 2018), and the CLIP vision-language model (Radford et al., 2021).

We compare our method to several competing approaches, including the class-agnostic Raw Attention, Rollout (Abnar & Zuidema, 2020), LRP (Bach et al., 2015), and Partial LRP (Voita et al., 2019) techniques, and the class-specific GradCAM (Selvaraju et al., 2017) and Transformer Attribution (commonly referred to as TransAttr) (Chefer et al., 2021b) methods. Although the various LRP method variants can generate explainable maps for each class, Chefer et al. (2021b) demonstrate that, in practice, the visualizations across different classes are largely similar. Thus, we treat these methods as class-agnostic. To reproduce the results of the competing methods, we follow the implementation provided by Chefer et al. (2021b). In all our experiments, we use standard ViT-B/16 and ViT-L/16 models pretrained on ImageNet (Russakovsky et al., 2015). We evaluate two variants of our PredicAtt method, constructed from the prediction maps of the last and second-to-last layers (PredicAtt_L and PredicAtt_{L-1}). Our implementation is based on the timm (Wightman, 2019).

Figure 5 shows qualitative comparisons on two representative images containing objects from different classes. As can be seen, our method better highlights the regions corresponding to each class.

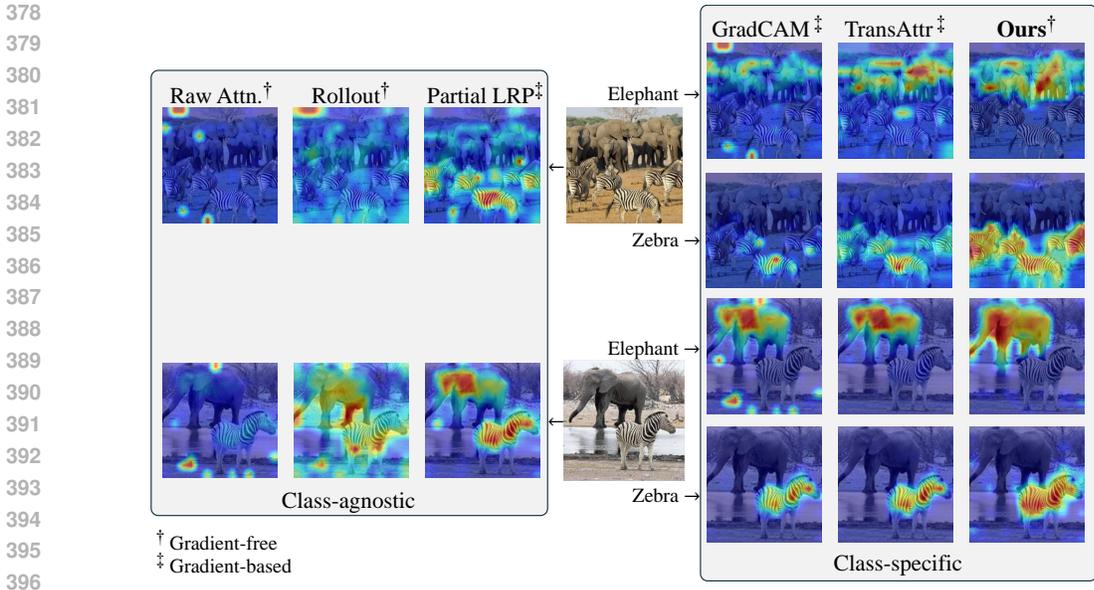


Figure 5: **Class-specific visualizations.** Our method, PredicAtt_{11} , captures a more coherent and compact region of the object. e.g. our method is the only one to highlight the elephant’s trunk.

We provide more visual examples in the supplementary. As suggested by Chefer et al. (2021b), we quantify the quality of the explanations using two measures, as follows:

Perturbation test. In this test, patches in the input image are gradually masked based on their importance, while measuring the model’s classification accuracy. This test has two variants; in the positive/negative version, pixels are masked in descending/ascending order of importance, leading to an expected sharp/gradual decline in accuracy. Both versions use the area under the curve (AUC) metric to quantify performance, scanning masking percentages of between 10% to 90% of the image. For the class-specific methods, we perform this experiment for both the predicted class and the target (ground-truth) class. Experiments are conducted on the ImageNet-Validation dataset (Russakovsky et al., 2015), which comprises 50,000 images from 1000 classes. The results are summarized in Tab. 1. As can be seen, both variants of our method achieve the best scores across all metrics.

		ViT-B/16				ViT-L/16			
		Negative		Positive		Negative		Positive	
Method		Pred. \uparrow	Target \uparrow	Pred. \downarrow	Target \downarrow	Pred. \uparrow	Target \uparrow	Pred. \downarrow	Target \downarrow
class-agnostic	Raw attention	45.55	-	24.00	-	40.91	-	27.22	-
	Rollout	53.10	-	20.06	-	52.75	-	21.67	-
	LRP	43.69	43.69	41.94	41.94	40.28	40.27	39.99	39.99
	Partial-LRP	50.29	50.28	19.81	19.82	37.23	37.24	29.56	29.56
class-specific	GradCAM	41.53	42.03	34.05	33.54	46.99	47.07	45.16	45.06
	TransAttr	54.20	55.09	17.04	16.41	51.75	52.40	20.03	19.61
	PredicAtt_{L-1} (Ours)	<u>55.41</u>	<u>56.99</u>	16.08	15.08	53.79	54.03	<u>19.98</u>	<u>19.78</u>
	PredicAtt_L (Ours)	56.16	57.94	<u>17.00</u>	<u>16.06</u>	<u>53.11</u>	<u>53.64</u>	19.10	18.66

Table 1: **Perturbation test.** All methods are evaluated on the ImageNet validation set with the ViT-B/16 and ViT-L/16 models. Bold and underline mark the best and second best scores, respectively. The subscript in our method indicates the layer of the prediction map, where L denotes the total number of layers in the model: 12 for ViT-B/16 and 24 for ViT-L/16.

Segmentation test. In this test, we use the explainability map generated from the predicted class to separate the foreground object from the background. To evaluate the performance of the segmentation, we employ three metrics: pixel accuracy, mean intersection-over-union (mIoU), and mean Average Precision (mAP), all computed based on the ground truth annotations. Our experiments are conducted on the ImageNet-Segmentation dataset (Guillaumin et al., 2014), comprising 4,276

		ViT-B/16			ViT-L/16		
Method		Pixel Acc.↑	mIoU↑	mAP↑	Pixel Acc.↑	mIoU↑	mAP↑
class-agnostic	Raw attention	67.87	46.37	80.24	63.20	41.18	74.75
	Rollout	73.54	55.42	84.76	71.15	52.88	83.48
	LRP	50.77	32.64	55.90	49.81	31.87	54.73
	Partial-LRP	76.31	57.97	84.67	62.40	40.21	73.65
class-specific	GradCAM	65.91	41.31	71.60	68.49	39.73	63.30
	TransAttr	<u>79.72</u>	<u>61.98</u>	86.04	72.88	52.20	81.22
	PredicAtt_{L-1} (Ours)	79.75	62.65	87.15	<u>78.32</u>	<u>59.20</u>	<u>84.38</u>
	PredicAtt_L (Ours)	76.85	59.22	<u>86.24</u>	83.06	64.51	86.78

Table 2: **Segmentation Test.** All methods are evaluated on the ImageNet-Segmentation dataset with the ViT-B/16 and ViT-L/16 models. Bold and underline mark the best and second best scores.

		ViT-B/16		ViT-L/16	
Method		Runtime ↓	Memory ↓	Runtime ↓	Memory ↓
class-agnostic	Raw attention	6 ms	8 MiB	11 ms	11 MiB
	Rollout	8 ms	28 MiB	19 ms	68 MiB
	LRP	158 ms	685 MiB	289 ms	2006 MiB
	Partial LRP	157 ms	685 MiB	285 ms	2006 MiB
class-specific	GradCAM	29 ms	511 MiB	42 ms	1637 MiB
	TransAttr	167 ms	681 MiB	341 ms	2013 MiB
	PredicAtt_L (Ours)	7 ms	68 MiB	20 ms	172 MiB

Table 3: **Resource consumption.** All results are for a single image. Our method is as fast and lightweight as gradient-free methods, running an order of magnitude faster and consuming an order of magnitude less memory than class-specific gradient-based methods.

images from 445 classes, each annotated with manual segmentation delineating the object in every image. Results are presented in Tab. 2. On ViT-L/16, both variants of our method achieve the highest scores across all metrics. On ViT-B/16, our PredicAtt_{L-1} variant achieves the highest scores on all metrics, while PredicAtt_L achieves the second-highest mAP score and third-highest pixel accuracy.

Resource consumption. Table 3 reports the GPU memory usage and runtime of all methods. Memory consumption refers to the difference between the peak memory allocation recorded post-generation of the explainable map and the memory allocated prior to this process. It thus quantifies the memory overhead of each method while neutralizing the memory footprint of the model weights. Notably, our method, which does not require gradients or backward passes, exhibits significantly shorter runtime and reduced memory consumption, highlighting its efficiency compared to competing methods. Experiments were conducted on an NVIDIA GeForce RTX 2080 Ti GPU.

6.1 ABLATION STUDY

We next study the contribution of each component in our approach by evaluating the following variants: (i) using only an average attention map, (ii) using only a weighted attention map, and (iii) using only a prediction map. The segmentation and perturbation tests for these variants with the ViT-B/16 model are reported in Tab. 4. While using only attention maps or only prediction maps yields unsatisfactory results, combining them leads to state-of-the-art performance.

Layer selection. Prediction maps can be obtained from any ViT layer. To determine which layers are most appropriate for explaining the prediction, we repeat the perturbation and segmentation tests for all layers. In Fig. 6, we observe a trend of improvement in all metrics as we take the prediction maps from deeper layers in the network. However, there is a slight degradation in the last layer, particularly in the segmentation and positive perturbation metrics. This may be attributed to the fact that the patch tokens at the output of the last layer do not affect the model’s output and are thus not

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

	Segmentation test			Perturbation test			
	Pix Acc↑	mAP↑	mIoU↑	Negative		Positive	
				Pred.↑	Target↑	Pred.↓	Target↓
(i) $\frac{1}{LH} \sum_{l,h} A_{CLS}^{(l,h)}$	67.38	80.73	44.52	49.45	49.45	23.08	23.08
(ii) $A_{CLS}(c)$	68.84	81.06	47.01	50.22	50.25	22.28	22.24
(iii) $\Psi^{(12)}(c)$	53.53	58.72	36.01	45.02	46.96	22.78	21.36
PredicAtt ₁₂ (c) (Ours)	76.85	86.24	59.22	56.16	57.94	17.00	16.06

Table 4: **Ablation study on ViT-B/16.** Using only attention maps (lines (i),(ii)) or only prediction maps (line (iii)), leads to weak explainability, yet their integration (Predicatt) yields state-of-the-art results. Note that a weighted average of the attention maps, which is based on their correlation with the prediction map (line (ii)), gives superior results to naively averaging them (line (i)).

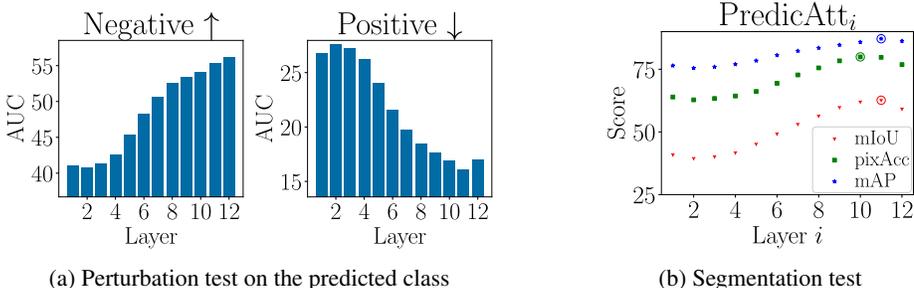


Figure 6: **Evaluation of PredicAtt_i per layer i.** The optimal score is marked with a circle. Performance generally improves with the use of deeper layers. However, in the segmentation and the positive perturbation tests, the prediction map from layer 11 outperforms that of the last (12th) layer.

optimized during training. Fig. 3 shows the prediction maps of the several last layers. The prediction maps of deeper layers tend to be more semantic, capturing the explained class more coherently.

7 CONCLUSION AND LIMITATIONS

We proposed a novel approach for explaining ViT predictions, which is based on a new visualization termed *Prediction Maps*. These maps, together with the well-studied attention maps form integral components of our explainability framework. We showed that the correlation between prediction- and attention-maps reliably indicates the influence of the attention maps on the predictions, and used this to construct a unified explainability map. Our method achieves state-of-the-art results in explainability measures and is significantly faster and more lightweight than the current leading methods. To the best of our knowledge, it is the first gradient-free method that provides class-specific explanations. In the supplementary, we discuss and illustrate the extension of our approach to explaining text models (BERT) and vision-language models (CLIP).

The primary limitation of our approach is its dependence on the classification head accepting tokens of a particular dimension. This assumption does not hold in several architectures, such as in DINO (Caron et al., 2021), where the classification head accepts a concatenation of [CLS] tokens from multiple layers, or in Swin (Liu et al., 2021b), where the size of the token embedding varies between layers. In these cases, adjustments to our method are necessary, and we leave them for future work.

REFERENCES

Samira Abnar and Willem Zuidema. Quantifying attention flow in transformers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4190–4197, 2020.

Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise

- 540 relevance propagation. *PLOS ONE*, 10(7):e0130140, 2015.
- 541
- 542 Homanga Bharadhwaj. Layer-wise relevance propagation for explainable deep learning based
543 speech recognition. In *2018 IEEE International symposium on signal processing and information
544 technology (ISSPIT)*, pp. 168–174. IEEE, 2018.
- 545 Adrien Bibal, Rémi Cardon, David Alfter, Rodrigo Wilkens, Xiaoou Wang, Thomas François, and
546 Patrick Watrin. Is attention explanation? an introduction to the debate. In *Proceedings of the 60th
547 Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.
548 3889–3900, 2022.
- 549 Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and
550 Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on
551 computer vision*, pp. 213–229. Springer, 2020.
- 552 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and
553 Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of
554 the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 9650–9660, October
555 2021.
- 556
- 557 Brandon Carter, Jonas Mueller, Siddhartha Jain, and David Gifford. What made you do this? under-
558 standing black-box decisions with sufficient input subsets. In *The 22nd International Conference
559 on Artificial Intelligence and Statistics*, pp. 567–576. PMLR, 2019.
- 560 Hila Chefer, Shir Gur, and Lior Wolf. Generic attention-model explainability for interpreting bi-
561 modal and encoder-decoder transformers. In *Proceedings of the IEEE/CVF International Confer-
562 ence on Computer Vision (ICCV)*, pp. 397–406, October 2021a.
- 563 Hila Chefer, Shir Gur, and Lior Wolf. Transformer interpretability beyond attention visualization.
564 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp.
565 782–791, 2021b.
- 566
- 567 Hila Chefer, Idan Schwartz, and Lior Wolf. Optimizing relevance maps of vision transformers
568 improves robustness. *Advances in Neural Information Processing Systems*, 35:33618–33632,
569 2022.
- 570
- 571 Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need
572 registers. *arXiv preprint arXiv:2309.16588*, 2023.
- 573 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
574 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 575
- 576 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
577 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An im-
578 age is worth 16x16 words: Transformers for image recognition at scale. In *International Confer-
579 ence on Learning Representations*, 2020.
- 580 Vijay Prakash Dwivedi and Xavier Bresson. A generalization of transformer networks to graphs.
581 *arXiv preprint arXiv:2012.09699*, 2020.
- 582
- 583 Ruth Fong, Mandela Patrick, and Andrea Vedaldi. Understanding deep networks via extremal per-
584 turbations and smooth masks. In *Proceedings of the IEEE/CVF international conference on com-
585 puter vision*, pp. 2950–2958, 2019.
- 586 Ruth C Fong and Andrea Vedaldi. Interpretable explanations of black boxes by meaningful pertur-
587 bation. In *Proceedings of the IEEE international conference on computer vision*, pp. 3429–3437,
588 2017.
- 589 Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel,
590 Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature
591 Machine Intelligence*, 2(11):665–673, 2020.
- 592
- 593 Matthieu Guillaumin, Daniel Küttel, and Vittorio Ferrari. Imagenet auto-annotation with segmenta-
tion propagation. *International Journal of Computer Vision*, 110:328–348, 2014.

- 594 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
595 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
596 770–778, 2016.
- 597 Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul
598 Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A criti-
599 cal analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF international
600 conference on computer vision*, pp. 8340–8349, 2021a.
- 602 Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial
603 examples. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
604 pp. 15257–15266. IEEE Computer Society, 2021b.
- 605 Sarthak Jain and Byron C Wallace. Attention is not explanation. In *Proceedings of the 2019 Con-
606 ference of the North American Chapter of the Association for Computational Linguistics: Human
607 Language Technologies, Volume 1 (Long and Short Papers)*, pp. 3543–3556, 2019.
- 609 Sonia Joseph. Vit prisma: A mechanistic interpretability library for vision transformers. <https://github.com/soniajoseph/vit-prisma>, 2023.
- 611 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
612 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceed-
613 ings of the IEEE/CVF International Conference on Computer Vision*, pp. 4015–4026, 2023.
- 614 Puneet Kumar, Vishesh Kaushik, and Balasubramanian Raman. Towards the explainability of mul-
615 timodal speech emotion recognition. In *Interspeech*, pp. 1748–1752, 2021.
- 617 Meir Yossef Levi and Guy Gilboa. Fast and simple explainability for point cloud networks. *arXiv
618 preprint arXiv:2403.07706*, 2024.
- 619 Jiwei Li, Xinlei Chen, Eduard Hovy, and Dan Jurafsky. Visualizing and understanding neural models
620 in nlp. *arXiv preprint arXiv:1506.01066*, 2015.
- 622 Naihan Li, Shujie Liu, Yanqing Liu, Sheng Zhao, and Ming Liu. Neural speech synthesis with
623 transformer network. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33,
624 pp. 6706–6713, 2019.
- 625 Yanghao Li, Chao-Yuan Wu, Haoqi Fan, Karttikeya Mangalam, Bo Xiong, Jitendra Malik, and
626 Christoph Feichtenhofer. Mvitv2: Improved multiscale vision transformers for classification and
627 detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recog-
628 nition*, pp. 4804–4814, 2022.
- 629 Shengzhong Liu, Franck Le, Supriyo Chakraborty, and Tarek Abdelzaher. On exploring attention-
630 based explanation for transformer models in text classification. In *2021 IEEE International Con-
631 ference on Big Data (Big Data)*, pp. 1193–1203. IEEE, 2021a.
- 633 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.
634 Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the
635 IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021b.
- 636 Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances
637 in neural information processing systems*, 30, 2017.
- 638 Girik Malik, Dakarai Crowder, and Ennio Mingolla. Extreme image transformations affect humans
639 and machines differently. *Biological Cybernetics*, 117(4):331–343, 2023.
- 641 Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey
642 Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Sim-
643 ple open-vocabulary object detection. In *European Conference on Computer Vision*, pp. 728–755.
644 Springer, 2022.
- 645 Ishan Misra, Rohit Girdhar, and Armand Joulin. An end-to-end transformer model for 3d object
646 detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp.
647 2906–2917, 2021.

- 648 Akash Kumar Mohankumar, Preksha Nema, Sharan Narasimhan, Mitesh M Khapra, Balaji Vasan
649 Srinivasan, and Balaraman Ravindran. Towards transparent and explainable attention models.
650 *arXiv preprint arXiv:2004.14243*, 2020.
- 651 Vitali Petsiuk, Abir Das, and Kate Saenko. RISE: randomized input sampling for explanation
652 of black-box models. In *British Machine Vision Conference 2018, BMVC 2018, Newcastle,
653 UK, September 3-6, 2018*, pp. 151. BMVA Press, 2018. URL [http://bmvc2018.org/
654 contents/papers/1064.pdf](http://bmvc2018.org/contents/papers/1064.pdf).
- 655 Zheng Qin, Hao Yu, Changjian Wang, Yulan Guo, Yuxing Peng, and Kai Xu. Geometric trans-
656 former for fast and robust point cloud registration. In *Proceedings of the IEEE/CVF conference
657 on computer vision and pattern recognition*, pp. 11143–11152, 2022.
- 658 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
659 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
660 models from natural language supervision. In *International conference on machine learning*, pp.
661 8748–8763. PMLR, 2021.
- 662 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. ”why should I trust you?”: Explaining the
663 predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference
664 on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, pp.
665 1135–1144, 2016.
- 666 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
667 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual
668 recognition challenge. *International journal of computer vision*, 115:211–252, 2015.
- 669 Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh,
670 and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based local-
671 ization. In *Proceedings of the IEEE international conference on computer vision*, pp. 618–626,
672 2017.
- 673 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In
674 *International conference on machine learning*, pp. 3319–3328. PMLR, 2017.
- 675 Gongbo Tang, Rico Sennrich, and Joakim Nivre. An analysis of attention mechanisms: The case
676 of word sense disambiguation in neural machine translation. *arXiv preprint arXiv:1810.07595*,
677 2018.
- 678 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
679 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
680 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 681 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
682 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-
683 tion processing systems*, 30, 2017.
- 684 Elena Voita, Pavel Serdyukov, Rico Sennrich, and Ivan Titov. Context-aware neural machine trans-
685 lation learns anaphora resolution. *arXiv preprint arXiv:1805.10163*, 2018.
- 686 Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head
687 self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *Proceedings
688 of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for
689 Computational Linguistics, 2019.
- 690 Sarah Wiegrefe and Yuval Pinter. Attention is not not explanation. In *2019 Conference on Empirical
691 Methods in Natural Language Processing and 9th International Joint Conference on Natural Lan-
692 guage Processing, EMNLP-IJCNLP 2019*, pp. 11–20. Association for Computational Linguistics,
693 2019.
- 694 Ross Wightman. Pytorch image models. [https://github.com/rwightman/
695 pytorch-image-models](https://github.com/rwightman/pytorch-image-models), 2019.

- 702 Junyi Wu, Bin Duan, Weitai Kang, Hao Tang, and Yan Yan. Token transformation matters: Towards
703 faithful post-hoc explanation for vision transformer. *arXiv preprint arXiv:2403.14552*, 2024.
704
- 705 Lian Xu, Wanli Ouyang, Mohammed Bennamoun, Farid Boussaid, and Dan Xu. Multi-class to-
706 ken transformer for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF*
707 *conference on computer vision and pattern recognition*, pp. 4310–4319, 2022.
- 708 Zhitao Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. Gnnexplainer:
709 Generating explanations for graph neural networks. *Advances in Neural Information Processing*
710 *Systems*, 32, 2019.
- 711 Hao Yuan, Haiyang Yu, Jie Wang, Kang Li, and Shuiwang Ji. On explainability of graph neural
712 networks via subgraph explorations. In *International conference on machine learning*, pp. 12241–
713 12252. PMLR, 2021.
- 714 Seongjun Yun, Minbyul Jeong, Raehyun Kim, Jaewoo Kang, and Hyunwoo J Kim. Graph trans-
715 former networks. *Advances in neural information processing systems*, 32, 2019.
- 716 Omar Zaidan, Jason Eisner, and Christine Piatko. Using “annotator rationales” to improve machine
717 learning for text categorization. In *Human language technologies 2007: The conference of the*
718 *North American chapter of the association for computational linguistics; proceedings of the main*
719 *conference*, pp. 260–267, 2007.
- 720 Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In
721 *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 16259–16268,
722 2021.
- 723 Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei
724 Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from
725 a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF confer-*
726 *ence on computer vision and pattern recognition*, pp. 6881–6890, 2021.
- 727 Tianhang Zheng, Changyou Chen, Junsong Yuan, Bo Li, and Kui Ren. Pointcloud saliency maps.
728 In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 1598–1606,
729 2019.
- 730 Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep
731 features for discriminative localization. In *Proceedings of the IEEE conference on computer*
732 *vision and pattern recognition*, pp. 2921–2929, 2016.
- 733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A SUPPLEMENTARY MATERIAL

A.1 PREDICTION MAPS PER LAYER

In this section, we provide quantitative and qualitative results of prediction maps from different layers of ViT-B/16 without combining them with attention maps. Figure 7 shows the perturbation and segmentation tests for the prediction map of each layer, while Fig. 8 shows visual results.

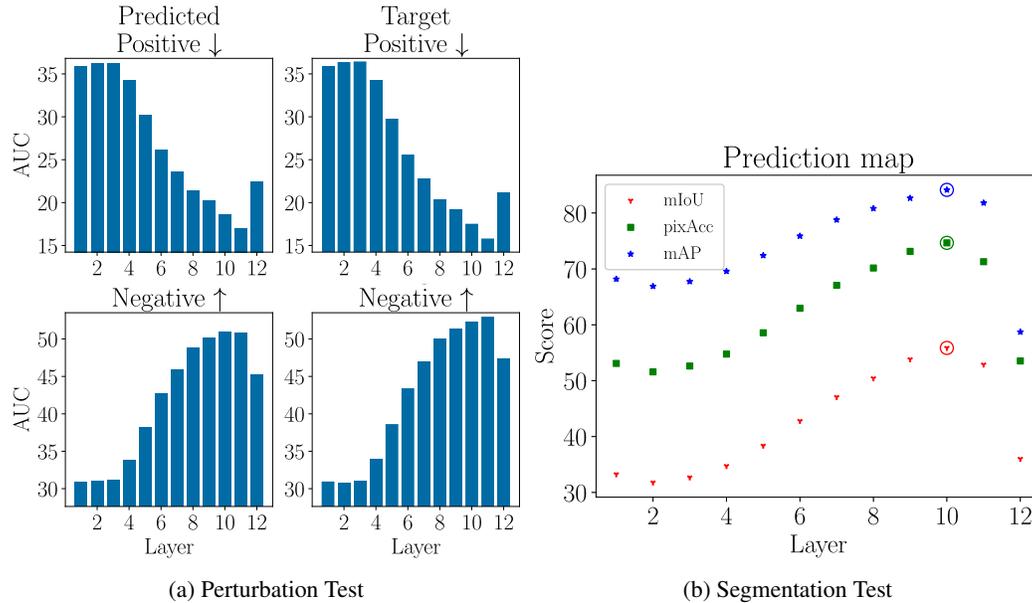
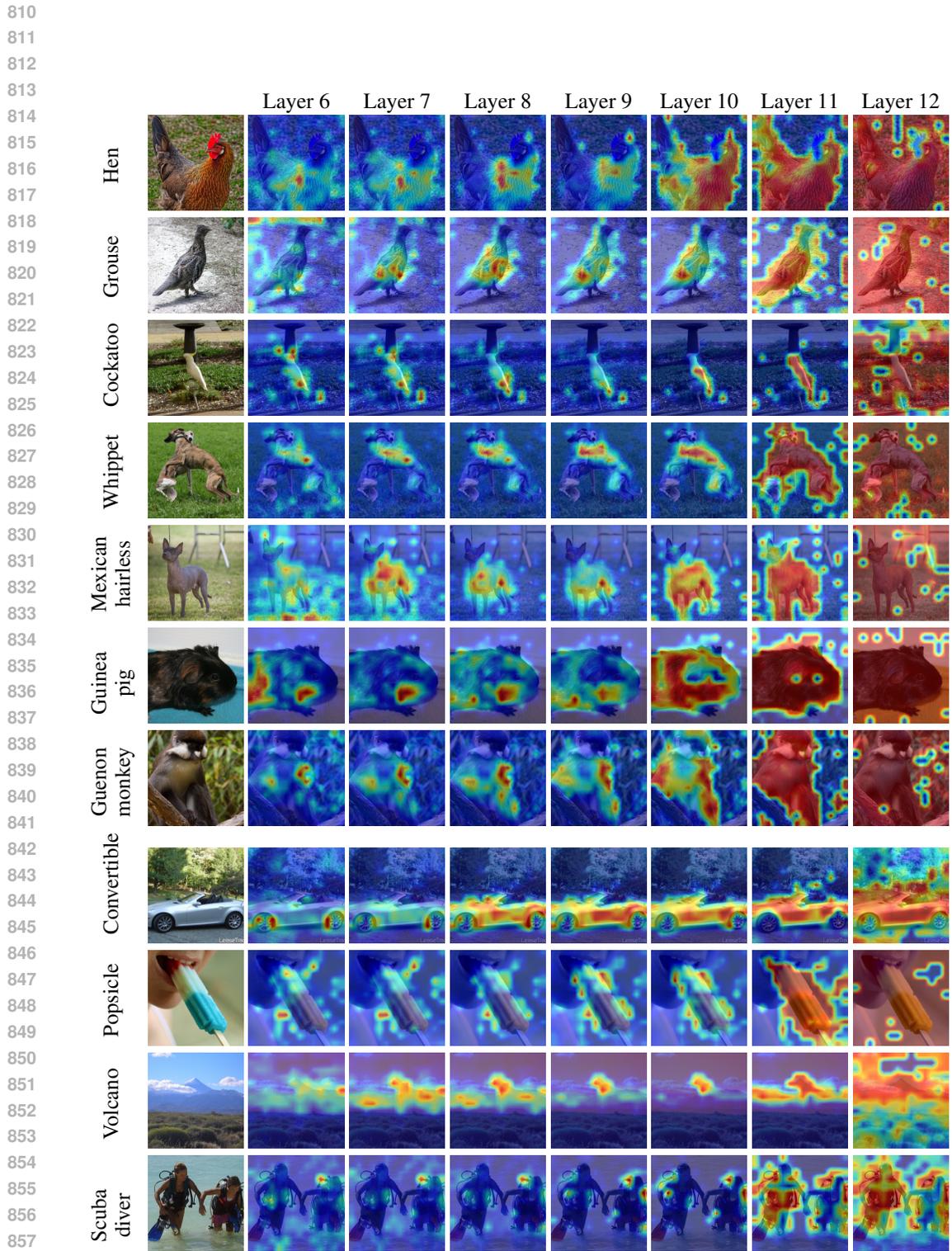


Figure 7: **Evaluation of prediction map per layer.** The optimal score is marked with a circle. The perturbation tests were run on a subset of the ImageNet validation dataset. Performance generally improves with the use of deeper layers. However, in perturbation tests the prediction map from layer 11 outperforms that of the last (12th) layer, and in the segmentation tests layer 10 outperforms layers 11 and 12.



859 **Figure 8: Prediction map per layer.** Additional examples on images from the ImageNet validation
860 dataset.

861
862
863

A.2 ALTERNATIVE SIMILARITY MEASURE

We now examine an alternative similarity measure for comparing the prediction map with attention maps, which forms the basis for generating the PredicAtt visualization. Specifically, instead of Eq. 4 we use the normalized dot product (correlation coefficient),

$$\alpha_{l,h} = \frac{\langle A_{\text{CLS}}^{(l,h)}, \Psi^{(i)}(c) \rangle}{\|A_{\text{CLS}}^{(l,h)}\| \cdot \|\Psi^{(i)}(c)\|}. \quad (6)$$

Table 5 compares the performance of PredicAtt₁₂ on the ViT-B/16 model using both the alternative and original similarity measures in segmentation tests. It is evident from the results that the alternative similarity measure is slightly inferior to the original.

Similarity Measure	Pixel Acc.↑	mIoU↑	mAP↑
Normalized dot-product	75.60	57.76	85.85
Dot-product	76.85	59.22	86.24

Table 5: **Comparison of Similarity Measures.** This table presents the performance results of PredicAtt₁₂ on the ViT-B/16 model during segmentation tests, evaluating both the original and alternative similarity measures. The results indicate that the alternative measure, the normalized dot product, demonstrates inferior performance compared to the original method.

A.3 ADDITIONAL VISUALIZATIONS OF PREDICATT

In Figures 9-15, we provide additional examples of the visualizations generated by our method, PredicAtt, as well as for the similarities between attention and prediction maps at different layers and heads.

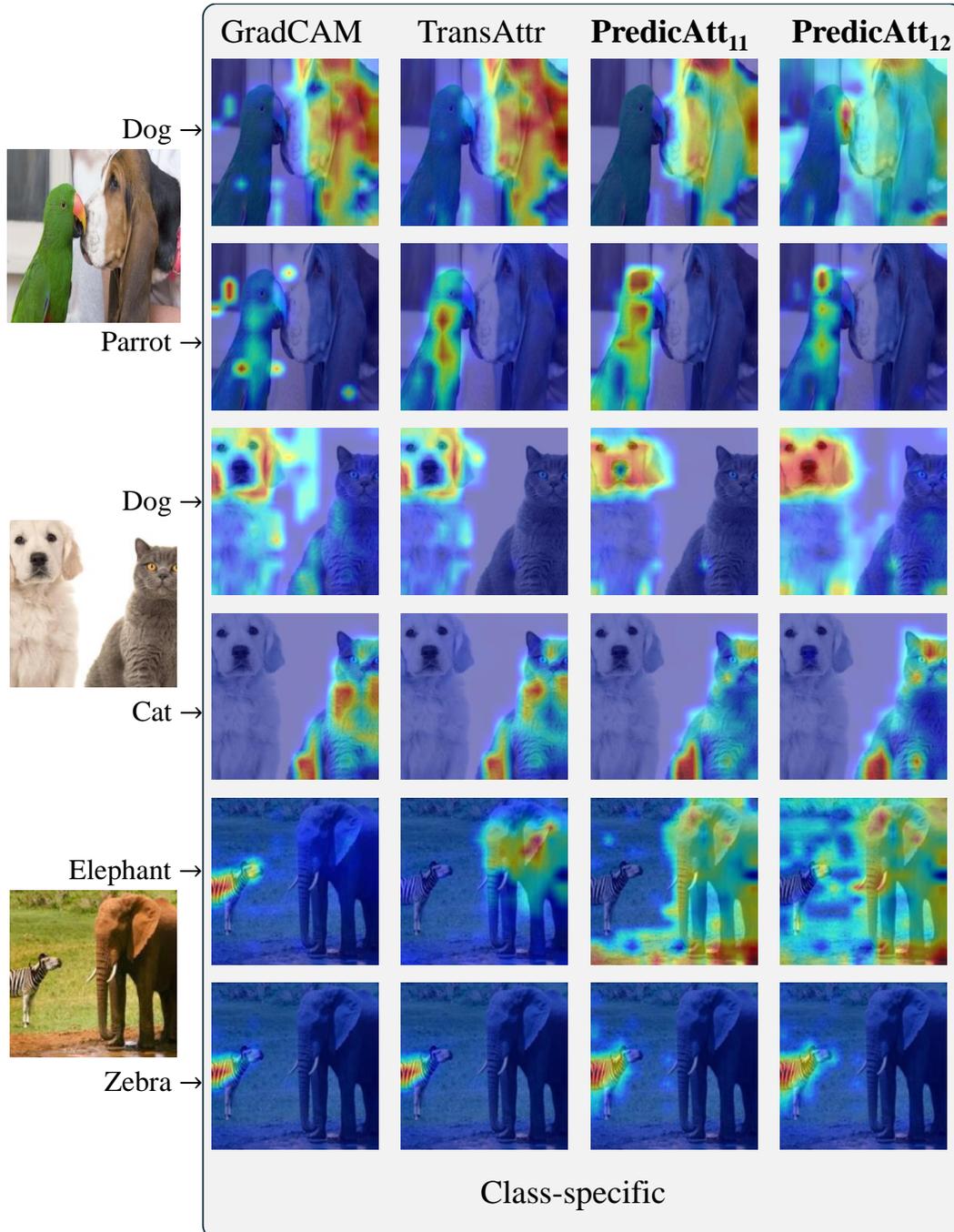


Figure 9: Class-specific visualizations.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

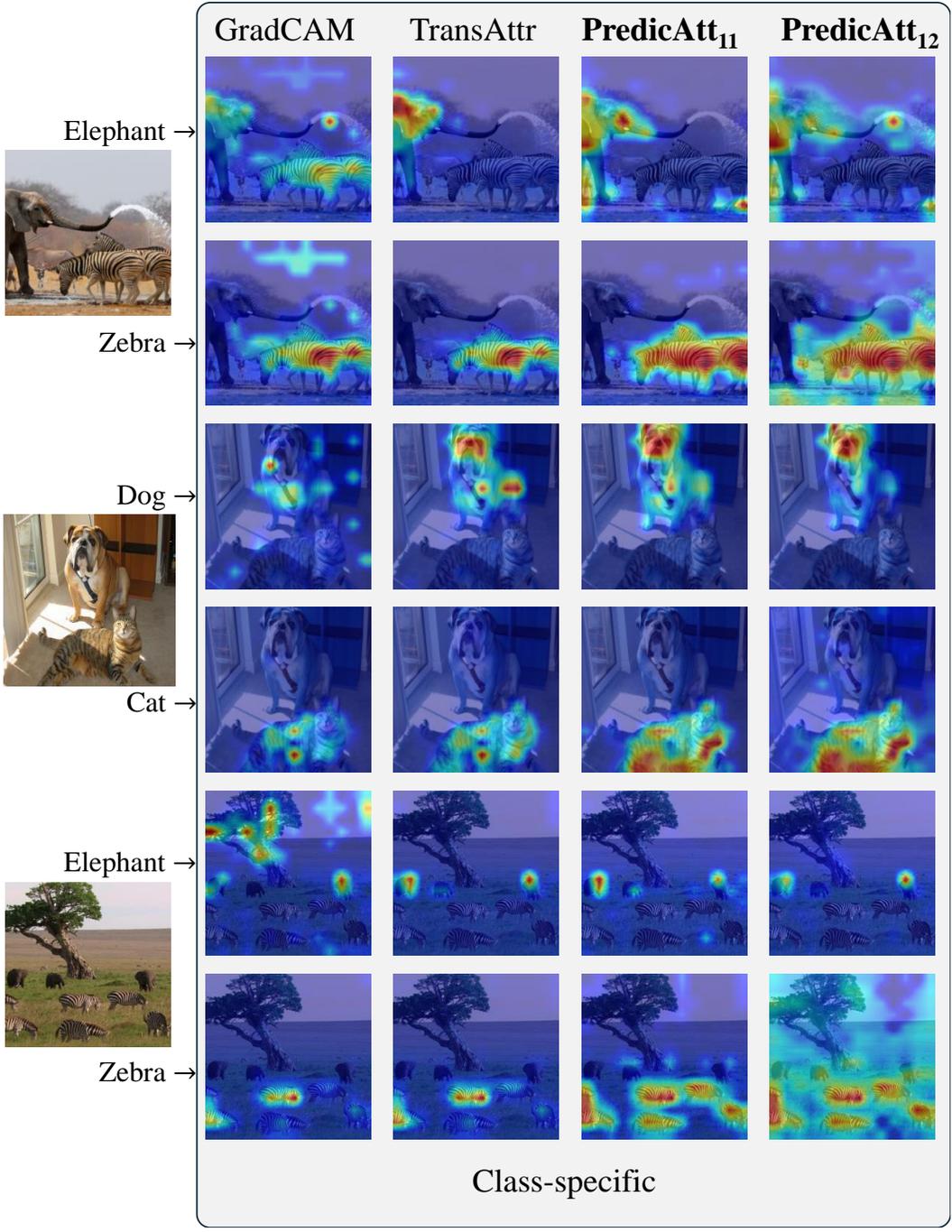


Figure 10: Class-specific visualizations.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

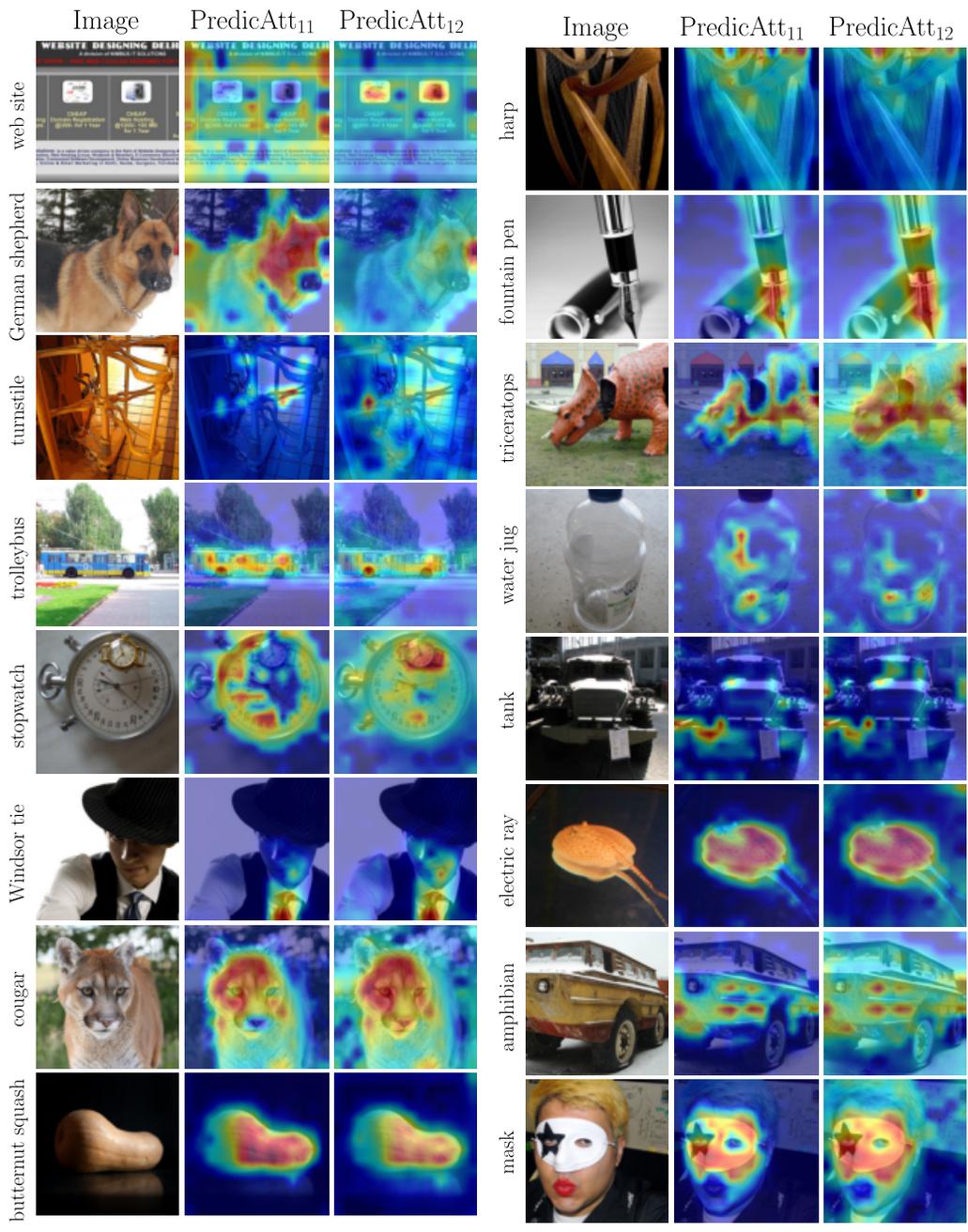


Figure 11: Sample images from ImageNet validation dataset. Some images are better explained by PredicAtt₁₁, while others by PredicAtt₁₂.

1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

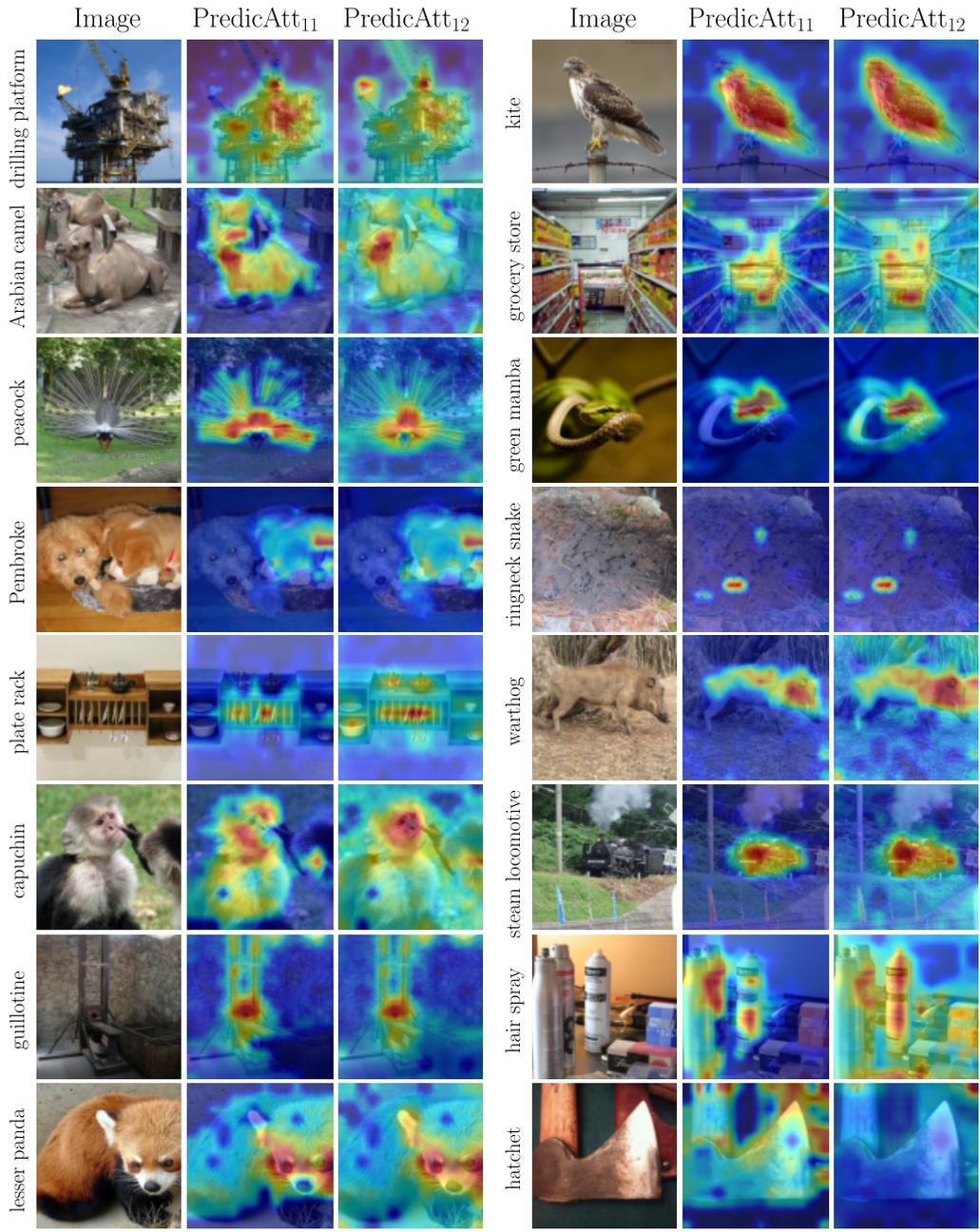


Figure 12: **Sample images from ImageNet validation set.** Some images are better explained by PredicAtt₁₁, while others by PredicAtt₁₂.

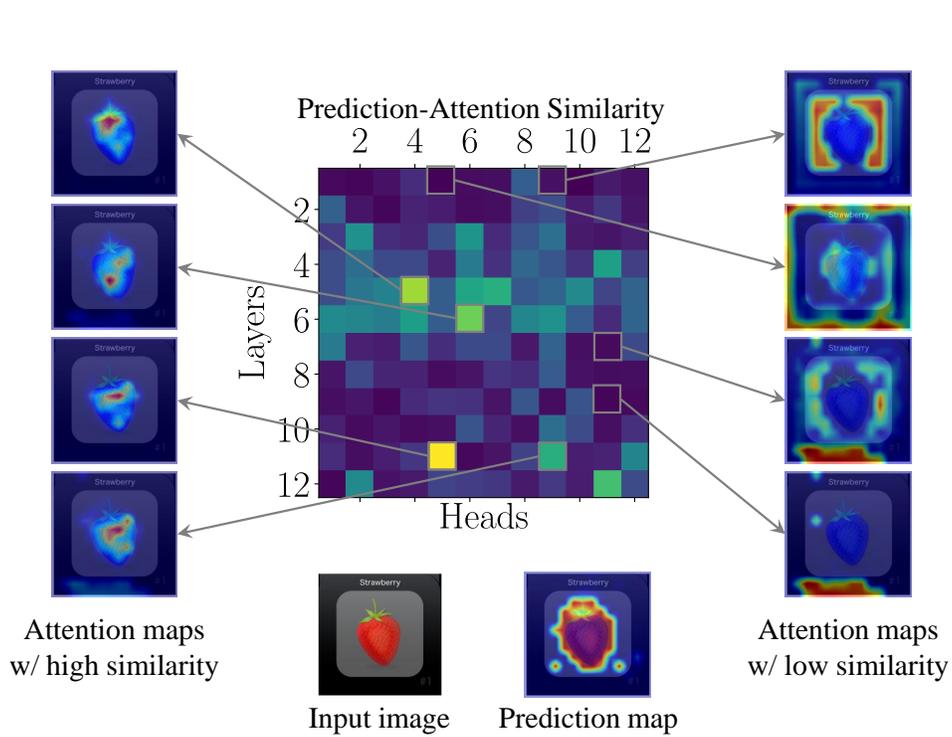


Figure 13: Per attention map analysis – Strawberry.

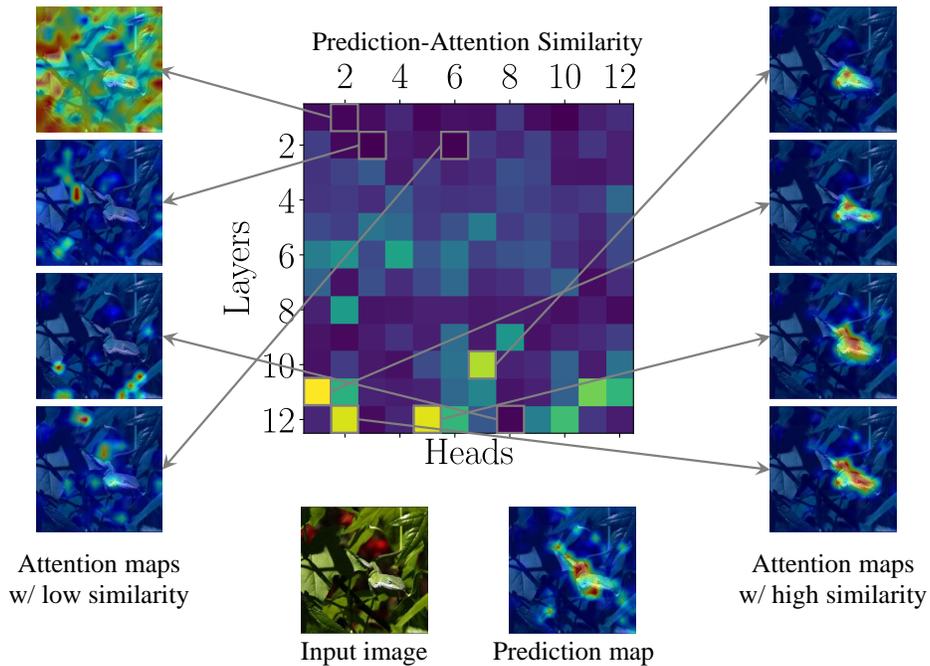


Figure 14: Per attention map analysis – Chameleon.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

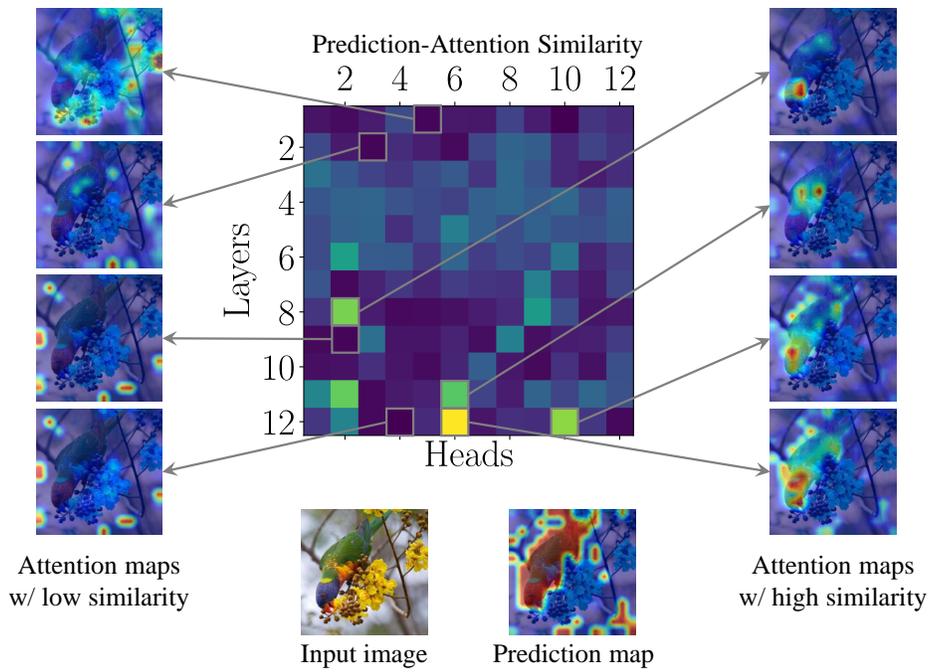


Figure 15: Per attention map analysis – Lorikeet.

1242 A.4 APPLYING PREDICTION MAPS TO TEXT

1243
1244 As discussed in Sec. 7, our method is applicable to any architecture with a classification head that can
1245 be fed with tokens other than the CLS token. We focused on image classifiers in order to compare
1246 to the existing explainability methods, which also focused on these models. But our method can
1247 be applied also to text classifiers with a similar architecture. In that modality, the explainable map
1248 would allow highlighting the word-parts in the text, which are most relevant for any desired class in
1249 a text-classification task.

1250 Figures 16 and 17 demonstrate the application of prediction maps to a BERT-base (Devlin et al.,
1251 2018) classifier model, assuming a maximum token sequence length of 512. Similarly to ViT, a
1252 classification token [CLS] is prepended to the input sequence and serves as the input to the classi-
1253 fication head. For this illustration, we use a BERT model fine-tuned on the Movie Reviews Dataset
1254 (Zaidan et al., 2007), a binary sentiment analysis task. The prediction maps visualize word parts
1255 most indicative of either positive or negative sentiment.

1256 this movie was the best movie i have ever seen ! some scenes were ridiculous but acting
1257 was great

1258 Figure 16: **Prediction Map for Positive Sentiment.** The prediction map highlights word segments
1259 associated with positive sentiment.

1262 i really didn't like this movie . some of the actors were good , but overall the movie was
1263 boring .

1264 Figure 17: **Prediction Map for Negative Sentiment.** The prediction map highlights word segments
1265 associated with negative sentiment.

1268 A.5 ADAPTING PREDICTION MAPS TO CLIP

1269
1270 In this section, we present an adaptation of our method to CLIP (Radford et al., 2021), a vision-
1271 language model. Specifically, given an image and a text, we would like to visualize how the CLIP
1272 image encoder associates between the text and each region in the image. The CLIP image encoder
1273 is not a classification model, and therefore it does not have a classification head. We propose to
1274 compute the cosine similarity between the text embedding obtained from the CLIP text encoder and
1275 the embedding of each patch token within a given layer in the CLIP image encoder (rather than doing
1276 so only for the final embedding at the output of the network). This yields a map that is similar in
1277 nature to the prediction maps we obtain for ViT classifiers. An illustration of this adaptation can be
1278 seen in Fig. 18. As opposed to our ViT classifier illustrations, here the heatmaps are inverted, with
1279 red indicating low similarity values and blue representing high similarity values. This is because we
1280 observed that tokens corresponding to the most relevant regions for the prompt rather tend to exhibit
1281 the lowest cosine similarity. This inversion occurs consistently across the Transformer layers. We
1282 leave the thorough analysis of this phenomenon to future research.

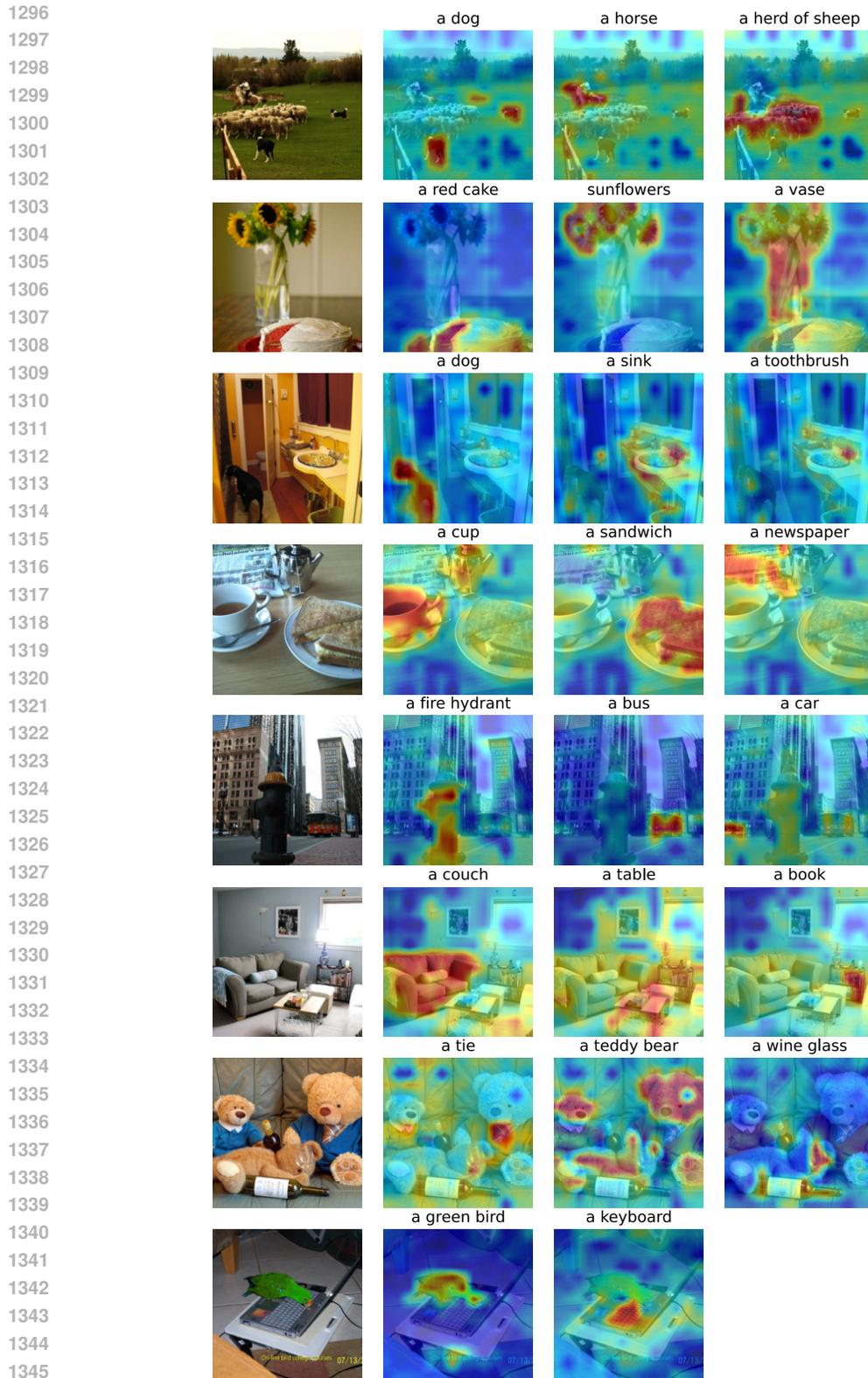


Figure 18: **CLIP explanation maps.** Adaptaion of our method for generating explainable heatmaps for the CLIP ViT-B/16 image encoder model. The patch tokens are extracted from the penultimate layer, with the caption above each heatmap indicating the corresponding prompt used for its generation.

A.6 ACCURACY, CONFIDENCE AND DISTRIBUTION OF THE PER-LAYER PREDICTIONS

In Figures 19,20,21, we analyze the accuracy, confidence and distribution of the predictions throughout the different layers. As a confidence measure, we utilize the common approach of computing the difference between the logits of the most probable class and the second most probable class, as done in Joseph (2023). As a distribution measure we report the mean entropy of the predictions per layer across the entire dataset. The analysis is based on a random batch of 80 images from the ImageNet validation set. The plots demonstrate the mean value of each metric, calculated separately for the CLS token and the patch tokens across all the sampled images.

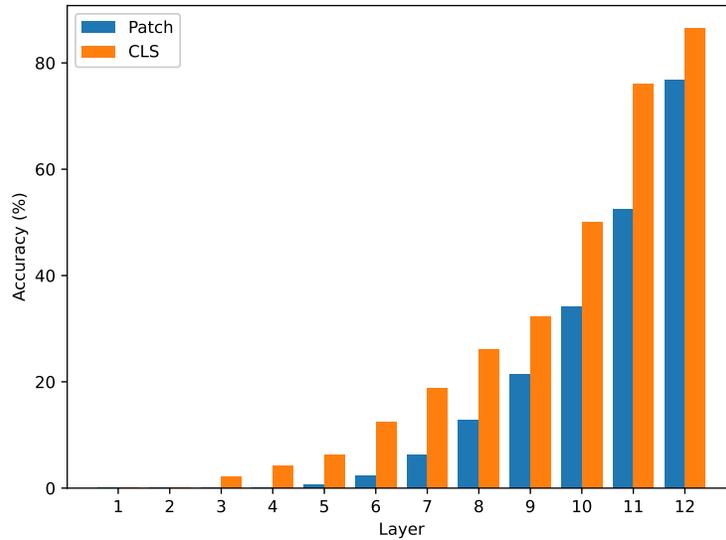


Figure 19: Accuracy per layer.

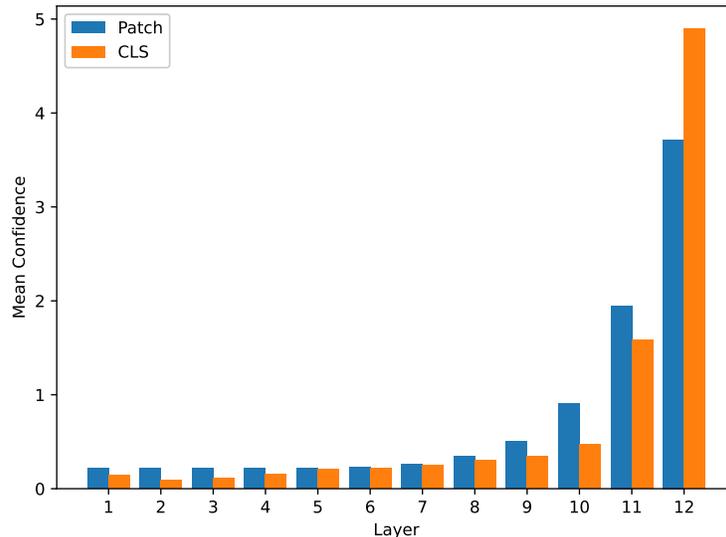


Figure 20: Mean confidence per layer.

1404
 1405
 1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414
 1415
 1416
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

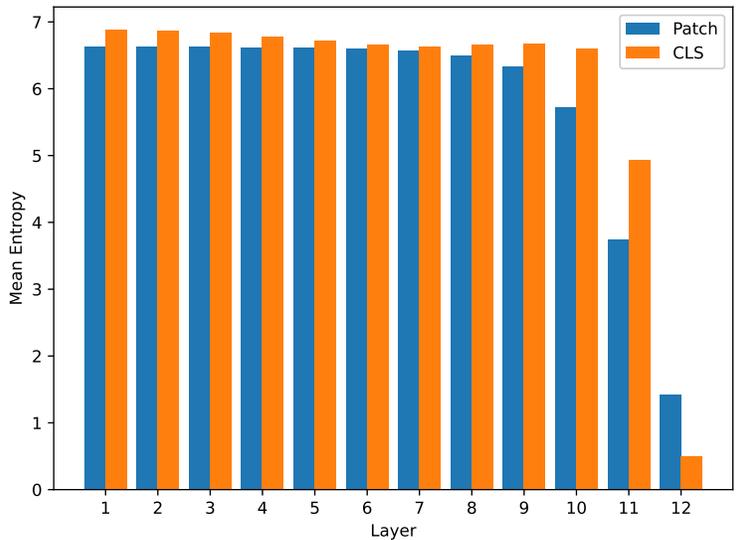


Figure 21: Mean entropy per layer.

We measure the similarity of the prediction map of the predicted class from the final layer to the attention heads across all layers. Figure 22 shows a histogram measuring the number of times the most similar attention head came from layer i , for every layer in the network. This analysis, conducted on the ImageNet-Segmentation dataset, reveals an unexpected pattern: while a monotonically increasing trend might be anticipated, the fifth layer emerges as the second most similar layer, surpassed only by the final layer.

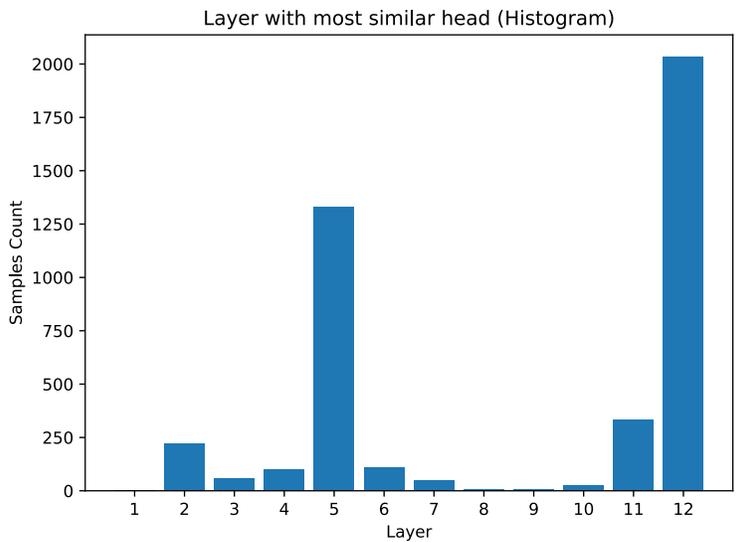


Figure 22: **Layer with the most correlated head.** Histogram of the layers whose attention heads are most similar to the final prediction map on ImageNet-Segmentation. The fifth layer shows notable correlation, second only to the final layer.

A.7 RESULTS ON SMALLER MODELS

In Tab. 6, we present the results of the perturbation and segmentation tests for ViT-S. Our method still demonstrates a slight improvement over TransAttr across most metrics.

Method	Perturbation test				Segmentation test		
	Negative		Positive		pixAcc \uparrow	mAP \uparrow	mIoU \uparrow
	Pred. \uparrow	Target \uparrow	Pred. \downarrow	Target \downarrow			
TransAttr	<u>53.22</u>	53.87	14.14	<u>13.78</u>	80.86	86.11	<u>63.61</u>
PredicAtt_{L-1} (Ours)	53.57	54.87	<u>14.24</u>	13.55	81.26	86.17	63.94
PredicAtt_L (Ours)	<u>53.22</u>	<u>54.61</u>	14.64	13.96	78.13	85.12	60.29

Table 6: **Perturbation and segmentation tests on ViT-S/16.** All methods are evaluated on the ImageNet validation set with the ViT-S/16 model. Bold and underline mark the best and second best scores, respectively. The subscript in our method indicates the layer of the prediction map, where L denotes the total number of layers in the model: 12 for ViT-S/16.