# Grad-ECLIP: Gradient-based Visual and Textual Explanations for CLIP

Chenyang Zhao, Kun Wang, Janet H. Hsiao and Antoni B. Chan

**Abstract**—Significant progress has been achieved on the improvement and downstream usages of the Contrastive Language-Image Pre-training (CLIP) vision-language model, while less attention is paid to the interpretation of CLIP. We propose a Gradient-based visual and textual Explanation method for CLIP (Grad-ECLIP), which interprets the matching result of CLIP for specific input image-text pair. By decomposing the architecture of the encoder and discovering the relationship between the matching similarity and intermediate spatial features, Grad-ECLIP produces effective heat maps that show the influence of image regions or words on the CLIP results. Different from the previous Transformer interpretation methods that focus on the utilization of self-attention maps, which are typically extremely sparse in CLIP, we produce high-quality visual explanations by applying channel and spatial weights on token features. Qualitative and quantitative evaluations verify the effectiveness and superiority of Grad-ECLIP compared with the state-of-the-art methods. Furthermore, a series of analysis are conducted based on our visual and textual explanation results, from which we explore the working mechanism of image-text matching, the strengths and limitations in attribution identification of CLIP, and the relationship between the concreteness/abstractness of a word and its usage in CLIP. Finally, based on the ability of explanation map that indicates text-specific saliency region of input image, we also propose an application with Grad-ECLIP, which is adopted to boost the fine-grained alignment in the CLIP fine-tuning. The code of Grad-ECLIP is available here: https://github.com/Cyang-Zhao/Grad-Eclip.

**Index Terms**—gradient-based explanation, visual and textual explanation, explainable AI, contrastive language-image pre-training, fine-grained understanding, open-vocabulary detection, deep learning

✦

## 1 INTRODUCTION

Recently, by learning the representations for matching caption text and its corresponding image, the Contrastive Language-Image Pre-training (CLIP) model [1] has introduced a simple and effective dual-encoder pre-training paradigm for the interaction between natural language processing and computer vision. CLIP significantly improves the performance on various downstream tasks, such as classification [2, 3], retrieval [4] and segmentation [5, 6], with zero-shot and fine-tuning methodologies. Inspired from CLIP, multi-modal pre-training has been further developed by exploring different perspectives, including unifying vision-language understanding and generation [7, 8], prompt design [9, 10], and region-aware enhancement [11, 12, 13]. Although researchers devote many efforts into improving multi-modal pre-training or exploring the usages in downstream tasks, less attention has been focused on the interpretation or explanation of CLIP.

Previous visual explanation works have considered interpreting the transformer architecture used by CLIP. Attention Rollout [14] generates explanations by aggregating attention maps computed along the forward pass of the model. Relevance-based methods [15, 16] apply Layer-wise Relevance Propagation (LRP) [17] and also rely on the attention mechanism in the model architecture. Since Rollout and many LRP variants are class-agnostic, Transformer interpretability [15] and Generic Attention-Model Explainability (GAME) [16] build class-specific relevance-based explanations using the self-attention or co-attention. However, just
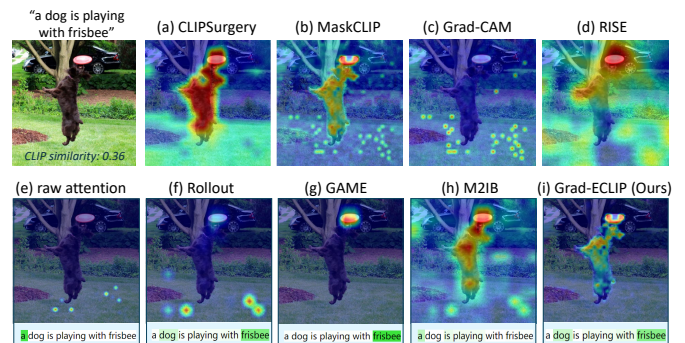


Fig. 1: Visual and textual explanations of CLIP for the image with the text "A dog is playing with frisbee" using (a) CLIPSurgery [18]; (b) MaskCLIP [19]; (c) Grad-CAM [20]; (d) RISE [21]; (e) raw attention in the last layer; (f) Rollout [14]; (g) GAME [16]; (h) M2IB [22]; and (i) Our Grad-ECLIP. For (e) to (i), textual explanations on the sentence are shown, where the degree of green color represents the word importance. Other methods (a-d) are not applicable on text.

treating CLIP as a vision transformer (ViT) and generating visual explanations based on self-attention sometimes leads to confusing results because of sparse attention maps (see Fig. 1e-g).

ECLIP [23] and CLIPSurgery [18] (see Fig. 1a) explores explanations for CLIP by computing an image-text similarity map, and solve the counter-intuitive problem that background patch features get higher similarity with the text feature than the foreground. However, to obtain reasonable similarity maps, new additional projection layers or changing the structure of the original CLIP are required. Although the parameters of CLIP encoders are frozen, learning more black-box parameters with extra data or modifying original model architecture makes the explanation less interpretable. MaskCLIP [19] also provides a technique to calculate class-specific image-text similarity map. By passing the

• *Chenyang Zhao, and Antoni B. Chan (corresponding author) are with the Department of Computer Science, City University of Hong Kong. Janet H. Hsiao is with the Division of Social Science and Department of Computer Science & Engineering, Hong Kong University of Science & Technology, and Kun Wang is with the SenseTime Group Ltd. E-mail: zhaocy2333@gmail.com, abchan@cityu.edu.hk.*

value features of the last attention layer through later linear layers as image patch features, the similarity map is able to localize the concept in the text (see Fig. 1b), but has noisy backgrounds and confusingly highlights points on the locations unrelated to the explained target. The disadvantage of these similarity-map methods is that they are only forward processing, and the attended features are not necessarily used in the final prediction.

To better focus on the discriminative features used in the prediction, gradient-based methods with class-activation maps (CAM) [24], such as Grad-CAM [20], Layer-CAM [25] and FullGrad [26], consider the gradient of the prediction with respect to features from a CNN layer as weights, and locates the class-specific discriminative regions by weighted aggregation of the features maps. Fig. 1c shows the visualization when adapting Grad-CAM on CLIP, where the cosine similarity of image-text pair is adopted as the prediction and the gradients are calculated w.r.t. the patch tokens from the ViT layers. Since there are no gradients w.r.t. the patch tokens in final layer because they are not involved in the calculation of the matching score, feature outputs from the penultimate layer of ViT are adopted. The results of Grad-CAM do not well explain CLIP, and suffer from the same problem of highlighting unrelated points as MaskCLIP, suggesting that the layer features of ViT are unsuitable for CAM methods.

In this paper, we explore a more effective way to interpret CLIP, by analyzing how CLIP obtains the final feature embedding, and deriving the relationship between the image and text embedding similarity score and intermediate features via a series of approximations. Based on the CAM principle, we then propose a novel gradient-based visual explanation method for CLIP (Grad-ECLIP), which generates the importance heat map by aggregating the intermediate image features with result-related *channel* and *spatial* importance weights. Our proposed method uses the gradients of the image-text matching score w.r.t. the attention layer as the importance for feature channels. For the spatial importance, because the softmax attention typically yields sparse attention maps, we propose a loosened attention map for computing the spatial importance, which can better reflect the importance of more regions, as compared to directly using the strict softmax attention. Then our Grad-ECLIP explanation map is calculated with the *values* in the attention layer as the feature map, weighted by the channel and spatial importances. The same method used to generate the explanation of the image encoder can also be applied to the text encoder to obtain a textual explanation for CLIP. Note that Grad-ECLIP is result-specific and is suitable for both the image and text encoders, i.e., the visual explanation on image is text-specific and the textual explanation of a sentence is image-specific. For the example shown in Fig. 1i, the heat map on the image shows the important region when matching the image with the specific text "a dog is playing with frisbee", while the degree of green color on the sentence represents the important words, where the most important words "dog" and "frisbee" correspond to the highlighted regions in the image heat map.

In the experiments, we conduct both qualitative evaluation by visualization of explanation maps and quantitative evaluation compared with other types of explanation methods, and show the superiority of our proposed Grad-ECLIP. Moreover, in the qualitative evaluation, we demonstrate the generalizability of Grad-ECLIP by applying our method on diverse datasets of different domains and showing it is applicable to both ViT and CNN-based CLIP, as well as the ViT classifier and other transformer-based vision language models like BLIP [8]. Then, using Grad-

ECLIP, we further conduct a visualization-based analysis on CLIP, and reveal working mechanisms and advantages/limitations of the CLIP model, including the type of attributes and the concreteness/abstractness of words used by CLIP. We hope our proposed method can be helpful for researchers to explore more properties of vision-language models like CLIP, as well as understand their current limitations and how this may affect downstream tasks.

Finally, we also present an application of Grad-ECLIP to fine-tuning the CLIP model to boost the fine-grained understanding. Since the ViT-based CLIP model has been shown to have limitations in producing dense representations, due to the pretraining focusing on the whole image-text matching [13, 12, 27, 28], we propose to generate detailed region-phrase corresponding pairs via Grad-ECLIP so as to enhance the fine-grained understanding of CLIP model during fine-tuning. Experiments with zero-shot region classification and down-stream open-vocabulary detection application show that the Grad-ECLIP-enabled fine-tuning is effective.

In summary, the contributions of this paper are five-fold:
1) We investigate Grad-ECLIP, a gradient-based visual and textual explanation approach for CLIP to produce high-quality result-specific heat maps for explaining the matching of image-text pairs.
2) We demonstrate the superiority of the proposed Grad-ECLIP with comprehensive evaluations comparing with the state-of-the-art explanation methods for Transformers and CLIP.
3) We show the generalizability of Grad-ECLIP by presenting explanations on datasets of different domains, and verifying the applicability on both ViT and CNN-based CLIP, ViT classifier, and other transformer-based vision language models.
4) By using Grad-ECLIP, we explore the properties of CLIP, and reveal the model's ability of concept decomposition and addibility, strengths and weaknesses in attribution identification, as well as the relationship between word usage and concreteness in the image and text matching.
5) We propose an application of Grad-ECLIP to boost the fine-grained understanding of CLIP via fine-tuning, which adopts the high-quality visual explanation map generated by Grad-ECLIP to produce detailed matching relationships between image regions and the corresponding semantic concepts.

A preliminary version of our work appears in [29]. The extensions over the conference version are as the follows. First, we provide a more detailed derivation of Grad-ECLIP, based on the transformer model. Second, we enrich the qualitative evaluation and verify the generalizability of Grad-ECLIP by adding: (1) the visualization on diverse datasets of different domains; (2) application on ViT-based classifier and other vision language models; (3) adaptation with CNN-based CLIP. Third, we include new ablation studies for the Grad-ECLIP design, including: (1) the effect of the proposed loosen spatial weight; (2) the influence of number of layers involved in the calculation; (3) the influence of multi-attention heads on the visual explanation results. Fourth, we add a new exploration about the relationship between word concreteness and word usage in image-text matching with CLIP. Moreover, we introduce an application of Grad-ECLIP, which applies the visual explanation map to boost the fine-grained region and text matching when fine-tuning CLIP.

The remainder of this paper is organized as follows. The related works about CLIP, explainability in computer vision and fine-grained understanding in CLIP are briefly reviewed in §2. Grad-ECLIP is introduced in §3, and the experiment results are presented in §4, including qualitative evaluations, quantitative

evaluations and ablation studies. We then perform analysis of CLIP based on the proposed Grad-ECLIP in §5. Finally, we introduce an application of Grad-ECLIP in §6 for boosting the fine-grained understanding of CLIP.

## 2 RELATED WORKS

We first briefly review the CLIP model, and then discus different types of visual explanation methods in computer vision. Finally, the related works for fine-grained image annotation and understanding in CLIP are introduced.

### 2.1 Contrastive language-image pre-training

Many multi-modal works have been developed and focus on the interaction of computer vision and natural language processing, such as text-image retrieval [30], image captioning [31], visual question answering [32], and visual grounding [33]. Contrastive language-image pre-training (CLIP) performs contrastive learning on very large-scale web-curated image-text pairs. It shows promising pre-trained representations with superior zero-shot transfer ability on diverse datasets and impressive fine-tuning performance on various downstream tasks. Subsequent works extend and improve CLIP from different aspects: [9, 10] improve the aspects of prompt design and optimization; [7, 8] unifies the vision-language understanding and generation by adding text decoders with image-text cross-attention during pre-training; [11, 12, 13, 28] builds an alignment between region feature or position information with fine-grained object descriptions. Although significant results have been achieved with CLIP and its development, less effort and exploration is focused on its interpretability through visual explanations. In this paper, we propose a novel visual explanation method, which generates high-quality heat maps that reveal the important regions or words used for CLIP's scoring of an image-text pair.

### 2.2 Explainability in computer vision

Since visualizing the importance of input features is a straightforward approach to interpret a model, many works visualize the internal representations of CNNs or Transformers with heat maps. Most explanation methods can be categorized as: CAM methods, perturbation methods, Shapley-value methods, or attribution propagation (relevance-based) methods.

CAM methods, such as CAM [24], Grad-CAM [20], and Grad-CAM++ [34], generate the explanation heat map from a selected feature layer by linearly aggregating the activation maps with weights that indicates each feature's importance. Grad-CAM computes the weights with global average pooling on the gradients of the model's prediction w.r.t. the feature layer. Gradient-free CAMs [35, 36, 37] generate weights from the prediction score changes when perturbing features.

Perturbation-based methods [38, 21, 39, 40, 41, 42, 43] perturb the input and observe the changes in output scores to determine the importance of input regions. Such black-box methods are intuitive and highly generalizable to different architecture and tasks, but computationally intensive. The quality of these methods are often greatly influenced by the type or resolution of the perturbations used. While having solid theoretical justification, Shapley-value methods [40] also suffer from large computational complexity.

Attribution propagation methods recursively decompose the network output into the contribution of early layers, based on the Deep Taylor Decomposition (DTD) [44]. LRP [17] and its variants [40, 45, 46] propagate relevance from the prediction to the input image based on DTD and generate class-agnostic explanations, while Contrastive-LRP [47] and SG-LRP [48] generate class-specific explanations.

These previous works are mainly proposed for interpreting CNN-based models. Due to the introduction of self-attention mechanisms in Transformers, recent works [49, 50, 51] have also looked at visual explanations for the Transformer architecture. [14] proposed an Attention flow and Rollout method, which is based on all attention maps in the forward process of model. Since Rollout is class-agnostic, Transformer interpretability [15] and GAME [16] build class-specific relevance-based method for explaining transformer with the internal attention mechanism. However, we found that the explanation methods relying on attention maps in Transformer cannot generate satisfactory results with CLIP, possibly because the sparse attention patterns on the $\mathrm{softmax}$ map. The recent M2IB [22] applies information bottleneck principle to CLIP, which develop an optimization objective to find the compressed representations for both image features and text features. However, a series of hyper-parameters are adopted during the optimization, which limits the generalization in practical applications.

Finally, existing approaches for explaining CLIP [23, 18, 19], which use the cosine similarity map between the image local features and the text features as the explanation map, have the disadvantage that they are only based on the *forward (bottom-up) process* and thus the attended features are not necessarily used in the final prediction. In contrast, we propose Grad-ECLIP as an effective approach to interpret CLIP, which highlight features that have largest influence on the prediction as measured by the gradient, which is a top-down process.

### 2.3 Fine-grained image understanding with CLIP

CLIP and its variants [9, 7, 8] exhibit strong representation capabilities and exceptional generalizability through learning general visual-language representations by pre-training on noisy large-scale datasets. Despite the great achievements, CLIP has shown lack of the fine-grained alignment between image regions and text [13, 12, 27, 28] due to its *image-level* training, which matches an image as a whole to a text description. Thus, the model is unable to generate precise representations of an image region for grounding textual concepts, which will limit the performance of CLIP on the downstream tasks that require region-aware ability. For example, in dense prediction tasks, e.g. object detection and segmentation, the CLIP model is usually utilized as a classifier [52, 53] or the teacher in distillation [54, 55] to process cropped object patches to obtain region features. Some works such as F-vlm [56], CORA [57] and FC-Clip [58] adopt the frozen CLIP model as backbone to produce spatial feature maps, but they all choose the CNN-based CLIP, which can preserve more position information than the vision transformer (ViT-based) architecture. However, due to the image-level training, CLIP models still lack fine-grained alignment and are poor at generating precise image region representations [13, 12, 27, 59].

To mitigate this issue, recent works enhance the fine-grained understanding ability of CLIP by leveraging region-text alignments during pre-training [60, 11, 61, 13, 59]. Since no region-text annotations are provided in the image-text pair training data, most of these methods need to generate image regions with the corresponding text tags using off-the-shelf methods. Some works utilize the annotations in visual grounding datasets [62, 59, 63]
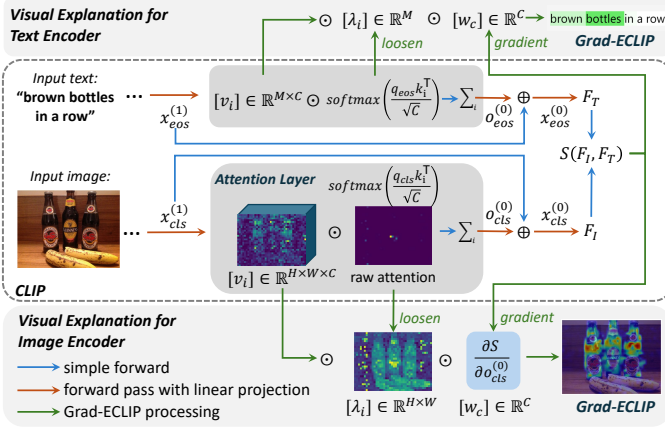
Fig. 2: Illustration of Grad-ECLIP. An image-text pair specific visual explanation is generated by weighting and aggregating the *values* as feature map in the attention layer with spatial importance $\lambda_i$ and channel importance $w_c$. Gradients are propagated to the attention layer output to produce $w_c$, and the loosened attention map is applied as $\lambda_i$.

or generate pseudo region-text pairs [60, 11, 61] with the help of high-performance detectors that are trained with a large number of object categories. RegionCLIP [13] adopts RPN [64] object proposals while PTP [12] coarsely crops patches, and then they both use CLIP as a classifier to obtain region labels with a large pre-defined pool of concepts, which are parsed from a text corpus. These methods inevitably cost significant extra time and space to preprocess the region annotations, which cannot be neglected when using a huge amount of training data. Moreover, the range of concepts is also limited by the number of pre-defined categories. Another work CLIPSelf [28] facilitates the transfer of the global features of the cropped regions to dense feature extraction by self-distillation, which enhances the local representations during the fine-tuning of CLIP model. However, preprocessing of generating region proposals by a well-trained detector is still required to obtain superior distillation performance – when using randomly cropped patches, the performance significantly drops.

In this paper, we build a novel fine-tuning framework for boosting the fine-grained understanding in CLIP, which adopts the proposed Grad-ECLIP to generate region-aware attention maps for aligning with the corresponding text phrases (concepts). By simply inserting the Grad-ECLIP-based module into the fine-tuning, our method circumvents the resource-consuming preparation of the region annotations and the requirement of high-performance detectors. The inputs of the proposed fine-tuning framework are still image-text pairs, just the same as in the pre-training of CLIP.

# 3 GRAD-ECLIP: GRADIENT-BASED VISUAL EXPLANATION FOR CLIP

Our method serves as a gradient-based visual and textual explanation for interpreting the CLIP matching performed on image-text pairs. We start with a brief introduction of CLIP. Then, by decomposing the layers of the transformer and exploring the relationship between the final output and intermediate features, we derive our gradient-based explanation for CLIP (Grad-ECLIP).

## 3.1 Preliminary on CLIP
CLIP learns both visual and language representations from large-scale web-curated image-text pairs. It consists of an image encoder $\mathcal{I}(\cdot)$ and a text encoder $\mathcal{T}(\cdot)$, which are jointly trained to respectively extract image and text feature embedding in a unified

representation space. Given image-text pair $(I, T)$, the matching score between their extracted image feature $F_I \in \mathbb{R}^D$ and text feature $F_T \in \mathbb{R}^D$ (both row vectors) is:

$$S(F_I, F_T) = \cos(F_I, F_T) = \frac{F_I F_T^{\mathsf{T}}}{\|F_I\| \|F_T\|}. \quad (1)$$

The model is trained using contrastive learning on the matching scores, regarding the ground-truth image-text pairs as positive samples and other mismatched pairs as negatives. In practice, both encoders can be implemented as transformers. In §3.2 our method is derived based on the transformer architecture, and thus is suitable for interpreting both ViT-based image and transformer-based text encoders. Alternatively, the image encoder could be a CNN-based ResNet followed by an attention pooling layer, which is basically the same as an attention layer in the Transformer. Thus the proposed Grad-ECLIP is also applicable to the CNN-based CLIP, for which we display the visualization results for CLIP with ResNet [65] backbone in §4.1.5.

## 3.2 Methodology of Grad-ECLIP
Here we present our derivation of Grad-ECLIP from the image viewpoint, where the visualization is generated on the input image $I$ and shows important regions related to producing the matching score $S_T(F_I) \triangleq S(F_I, F_T)$, with the given specific text prompt $T$. The application of Grad-ECLIP from the text viewpoint, where the visualization is generated for the text prompt $T$ given the input image $I$, can be obtained analogously by considering the $[eos]$ token (end of sentence token) from the text encoder, which is analogous to the $[cls]$ token in the image encoder.

For a Transformer that consists of $N$ layers, following convention, we denote $x^{(n)}$ as the input of layer $L^{(n)}$ and output of layer $L^{(n+1)}$, where $n \in [0...N]$ is the layer index. $x^{(N)}$ is the input of the network, $x^{(1)}$ is the input of the last layer and $x^{(0)} = \mathcal{I}(x^{(N)})$ is the output of the network. The image feature is $F_I = \mathcal{LP}(x_{cls}^{(0)})$, where $\mathcal{LP}$ denotes linear projections, and $x_{cls}$ is the feature vector from the $[cls]$ token. Thus, except for the class token, all the final layer features of the other tokens (image patch tokens) are not used during contrastive learning of CLIP. Therefore, to interpret the $S_T(F_I)$ w.r.t. an image feature map, we explore the relationship between the last layer class token feature $x_{cls}^{(0)} \in \mathbb{R}^C$ and the intermediate spatial feature maps.

As shown in the illustration of Fig. 2, looking closely into the last layer of the network, the image embedding from visual encoder can be formulated as:

$$F_I = \mathcal{LP}(x_{cls}^{(0)}) = \mathcal{LP}(o_{cls}^{(0)} + x_{cls}^{(1)}) \quad (2)$$
$$\approx \mathcal{LP}(o_{cls}^{(0)}) + \mathcal{LP}(x_{cls}^{(1)}), \quad (3)$$

where the approximation is based on assuming linearity of $\mathcal{LP}$. Noting that $\mathcal{LP}(x_{cls}^{(t)}) = \mathcal{LP}(o_{cls}^{(t)} + x_{cls}^{(t+1)})$, we substitute recursively to obtain the approximation:

$$F_I \approx \mathcal{LP}(o_{cls}^{(0)}) + \cdots + \mathcal{LP}(o_{cls}^{(N-1)}) + \mathcal{LP}(x_{cls}^{(N)}) \quad (4)$$
$$\triangleq \sum_{t=0}^{N} F_I^{(t)}, \quad (5)$$

where we define $F_I^{(t)} = \mathcal{LP}(o_{cls}^{(t)})$ for $t < N$, and $F_I^{(N)} = \mathcal{LP}(x_{cls}^{(N)})$. Thus from (5), the image feature $F_I$ is approximately an aggregation of features from each layer $F_I^{(t)}$.

The feature vector from each layer is computed from a self-attention operation. For example, for the last layer ($t = 0$),

$$F_I^{(0)} = \mathcal{LP}(o_{cls}^{(0)}) = \mathcal{LP}\big(\sum_i \text{softmax}(\frac{q_{cls}k_i^\mathsf{T}}{\sqrt{C}})v_i\big), \quad (6)$$

where the output of attention layer ($\mathcal{A}$) on the class token is

$$o_{cls}^{(0)} = \mathcal{A}(x^{(1)})[cls] = \sum_i \text{softmax}(\frac{q_{cls}k_i^\mathsf{T}}{\sqrt{C}})v_i, \quad (7)$$

and $\mathcal{A}$ represents the attention layer in the Transformer, $q_{cls}$ is the *query* embedding for the class token, while $k_i$ and $v_i \in \mathbb{R}^C$ represent the *key* and *value* embeddings at spatial location $i$, with $C$ as their channel dimension.[1] The softmax operation inside the attention layer measures the weight of the value on each location. Multi-heads are usually used in the attention layer to group the channel of $\{q, k, v\}$ into several heads, and (7) is operated inside each head with the softmax calculated over subsets of the channels. Then the final attention layer output is obtained by concatenating the results of each head together. In practice for visualization, we formulate the $o_{cls}^{(0)}$ with one attention head in the forward pass and operate the softmax over all channels as in (7). We discuss the influence multi-heads to visual explanation in the §4.3.3.

Assuming that the feature vectors $(F_T, F_I)$ are normalized and using (5), the matching score can be approximated as an aggregation of partial scores for feature vectors from each layer,

$$S_T(F_I) = \sum_c F_T[c]F_I[c] \approx \sum_c F_T[c]\sum_t F_I^{(t)}[c] \quad (8)$$

$$= \sum_t \Big(\sum_c F_T[c]F_I^{(t)}[c]\Big) \triangleq \sum_t S_T(F_I^{(t)}), \quad (9)$$

where $[c]$ selects the $c$-th channel, and $S_T(F_I^{(t)})$ denotes the score from layer $t$. Next, to calculate the heat map for the contribution of $o_{cls}$ on the partial matching score, we write the partial matching score as a function of its $o_{cls}$. Specifically, looking at the last layer ($t = 0$) as an example,

$$S_T(F_I^{(0)}) = \sum_c F_T[c]F_I^{(0)}[c] \quad (10)$$

$$= \sum_c F_T[c]\mathcal{LP}(o_{cls})[c]^{(0)} \triangleq f(o_{cls}), \quad (11)$$

where we have defined the matching score as a function of $o_{cls}$, i.e., $f(o_{cls})$. We define the approximation of the matching score as a weighted combination of the channel features in $o_{cls}$, we have

$$f(o_{cls}) \approx \tilde{f}(o_{cls}) \triangleq \sum_c w_c o_{cls}[c] = wo_{cls}^\mathsf{T} \quad (12)$$

where $w_c$ is the weight for the $c$-th channel, and $w = [w_c]_c \in \mathbb{R}^C$ the corresponding weight vector for all channels. To obtain the channel weights $w$, we aim to match the first derivatives (gradients) of the original matching score $f$ and its approximation $\tilde{f}$, leading to the optimization problem:

$$w = \underset{w}{\text{argmin}} \big\| f'(o_{cls}) - \tilde{f}'(o_{cls}) \big\|^2 \quad (13)$$

$$= \text{argmin} \big\| \frac{\partial f}{\partial o_{cls}} - w \big\|^2, \quad (14)$$

which has the solution

$$w = \frac{\partial f}{\partial o_{cls}} = \frac{\partial S_T(F_I)}{\partial o_{cls}}. \quad (15)$$

1. We skip the superscript (0) on $\{q, k, v\}$ identifying the layer for brevity.

Therefore, substituting (7) and (15) into (12),

$$S_T(F_I) \approx \sum_c w_c o_{cls}[c] \quad (16)$$

$$= \sum_c \frac{\partial S_T(F_I)}{\partial o_{cls}[c]} \sum_i \text{softmax}(\frac{q_{cls}k_i^\mathsf{T}}{\sqrt{C}})v_{ic} \quad (17)$$

$$= \sum_i \Big[\sum_c \frac{\partial S_T(F_I)}{\partial o_{cls}[c]}\text{softmax}(\frac{q_{cls}k_i^\mathsf{T}}{\sqrt{C}})v_{ic}\Big], \quad (18)$$

where $v_{ic}$ is the c-th channel of $v_i$. Regarding $[v_{ic}]$ as the intermediate feature map, and $\lambda_i = \text{softmax}(\frac{q_{cls}k_i^\mathsf{T}}{\sqrt{C}})$, then the importance of the $i$-th spatial location is defined as:

$$H_i = \text{ReLU}\big(\sum_c w_c\lambda_i v_{ic}\big) \quad (19)$$

where ReLU means that we only focus on positions that have a positive effect on the final score.

The weight $w_c$ represents the importance of each feature channel, and $\lambda_i$ represents the importance of the values at each location via the softmax attention. However, from the visualization, we discover that the output of the softmax self attention function is extremely sparse. Important information may be encoded in different locations, but the softmax only selects the largest activation, which is not appropriate as a spatial weight. Therefore, we replace the softmax with a "loosened" correlation by applying 0-1 normalization on the similarities $[q_{cls}k_i^\mathsf{T}]_i$, i.e., $\lambda_i \approx \Phi(q_{cls}k_i^\mathsf{T})$, where $\Phi$ is the 0-1 normalization function applied over the set of similarities. In the experiments and §4.3.1, we compare using the loosened $\lambda_i$ and without $\lambda_i$ to show the effect of spatial weights, qualitatively and quantitatively.

Therefore, with *the channel importance* $w_c$ and *the spatial importance* $\lambda_i$, where

$$w_c = \frac{\partial S_T(F_I)}{\partial o_{cls}^{(0)}[c]}, \quad \lambda_i = \Phi(q_{cls}k_i^\mathsf{T}), \quad (20)$$

we obtain the proposed Grad-ECLIP explanation map $H = [H_i]_i$ for the last layer, where $H_i$ is defined in (19) using the last layer values $v$ as the feature map.

**Visual and textual explanations:** For the image encoder, the flattened heat map $H_i$ is reshaped and interpolated to the original image's height and width, while for text encoder, the heat (importance) value on the $i_{th}$ tokens is remapped to the original word position in the sentence. Finally, based on the approximation in (9), the final explanation can be *aggregated over all the layers* by recursively processing each layer to obtain its heat map from (19). In the experiments, we use the last layer to explain the image encoder, and the last eight layers for interpreting the text encoder. The ablation study for the influence of different number of layers involved in image and text explanation is shown in §4.3.2.

**CNN-based CLIP:** Our proposed Grad-ECLIP can also be applied to CNN-based CLIP. The CNN-based CLIP model is composed of a ResNet backbone followed by an attention pooling. Thus, we can use the final attention layer in the pooling to conduct our explanation, which uses the same implementation as for ViT-based CLIP. The visualizations shown in §4.1.5 verify the effectiveness of Grad-ECLIP on CNN-based CLIP model.

# 4 EXPERIMENTS WITH GRAD-ECLIP

In this section we conduct experiments on Grad-ECLIP to: 1) evaluate its visual explanation qualitatively and quantitatively, and compare with the current SOTA methods; 2) evaluate the
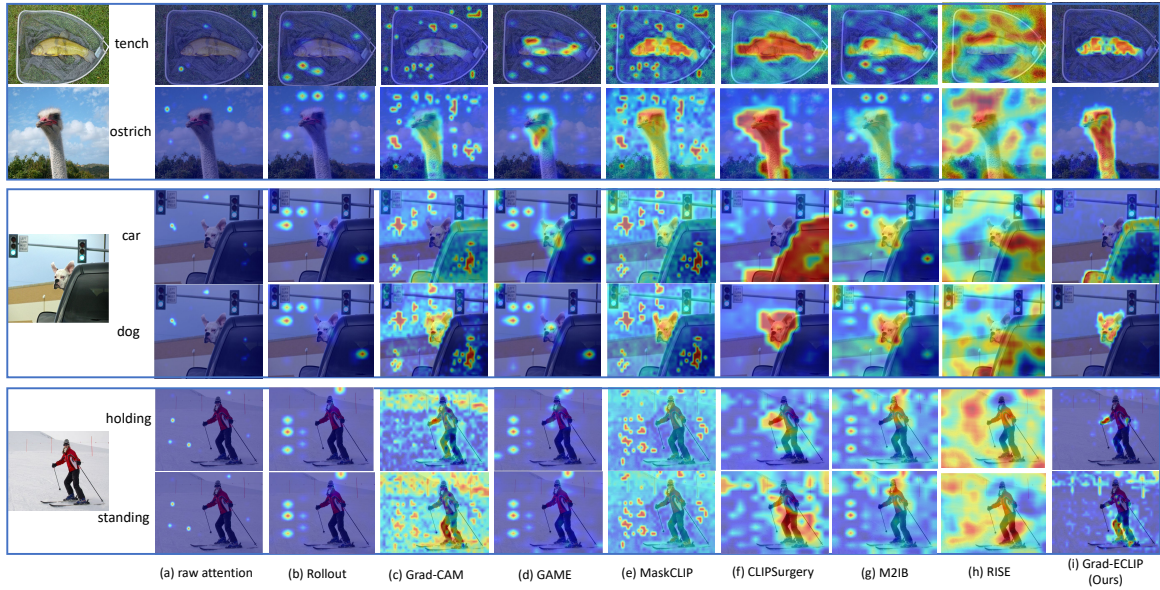
Fig. 3: Comparison of heat maps from: (a) the raw self-attention map in the last ViT layer; (b) Rollout [14]; (c) Grad-CAM [20]; (d) GAME [16]; (e) MaskCLIP [19]; (f) CLIPSurgery [18]; (g) M2IB [22]; (h) RISE [21]; (i) our proposed Grad-ECLIP. Visual explanations are provided for the matching score between the image and the specific text prompts, which can be nouns (*e.g.*, car, dog) or verbs (*e.g.*, holding, standing). From the comparison of visualizations, Grad-ECLIP exhibits superior explanation ability on different types of text prompts.



Fig. 4: Explanations for image-text pairs from MS COCO using: by (a) raw self-attention; Transformer interpretation methods (b) Rollout, (c) GAME; and our method (d) Grad-ECLIP. The importances of words are visualized by the degree of green color.

processing time; 3) conduct ablation studies, including the effect of spatial weight, the involved layers and attention heads.

Unless otherwise specified, we conducted the experiments with the ViT-B/16 architecture. We compared with representative baseline XAI methods from each category: 1) attention map-based *Rollout* [14], which takes into account all the attention maps computed along the forward pass, and *raw attention* in the last visual encoder layer, both of which are not result-specific explanation; 2) classical gradient-based method *Grad-CAM* [20], which takes the image-text similarity as target and calculate the gradients w.r.t. the ViT layer output; 3) relevance-based *GAME* [16], which

integrates the relevancies and gradients propagated through the network; 4) cosine-based *MaskCLIP* [19] and *CLIPSurgery* [18], which generates similarity value on each location by the cosine between text feature and processed values as local image features; 5) M2IB [22], which applies information bottleneck principle to generate explanation map for CLIP. Each baseline is built with different properties and assumptions over the architecture. We also show visualization comparisons with the typical black-box perturbation method RISE [21], but did not conduct quantitative comparisons with black-box perturbation and Shapely methods, due to their computational complexity and inherent differences
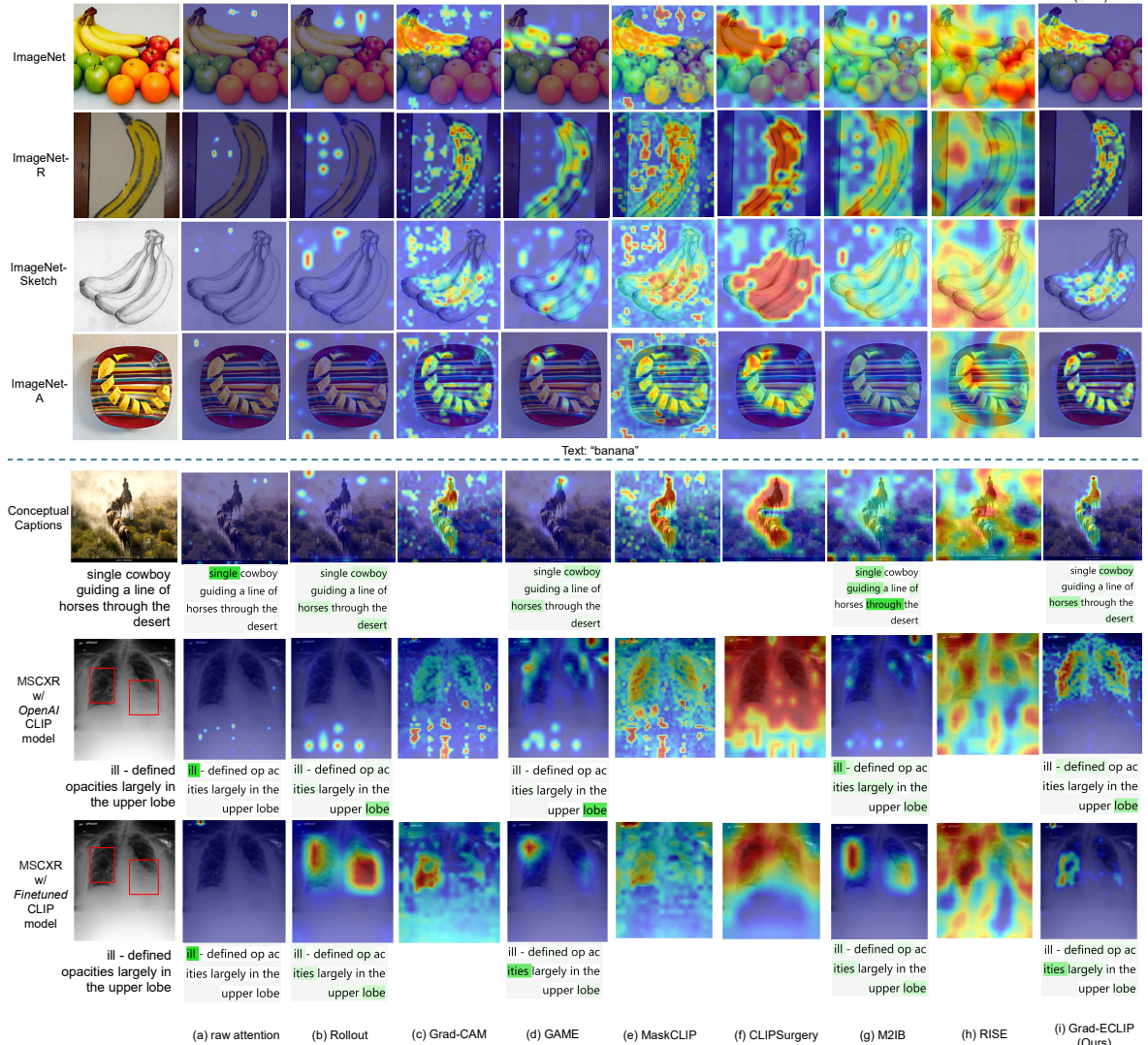
Fig. 5: Comparison of visual explanations for different methods on image samples from different image domains: natural images (ImageNet), renditions (ImageNet-R), pencil sketch (ImageNet-Sketch), natural adversarial examples (ImageNet-A), web images with captions (Conceptual Captions), and chest X-ray (MSCXR). On MSCXR, explanations are provided for both the OpenAI CLIP model and a fine-tuned version.

with our proposed approach.

## 4.1 Qualitative evaluation

In this section, we evaluate the proposed Grad-ECLIP qualitatively. The comparisons of visual explanations from different types of methods are presented in §4.1.1, and the visual explanations on both image and text encoder with image-sentence pair examples are preseneted in §4.1.2. Then, we conduct the visualizations of explanations on diverse image domains in §4.1.3 and adapt the Grad-ECLIP to other Vision Language Models (VLMs) in §4.1.4. Finally, we show that the Grad-ECLIP is applicable to CNN-based CLIP with ResNet backbone in §4.1.5.

### 4.1.1 Comparison of visual explanations

We compare the visualizations of raw self-attention, Rollout, Grad-CAM, GAME, MaskCLIP, CLIPSurgery, M2IB, RISE and our Grad-ECLIP in Fig. 3 with images from ImageNet [66] and MS COCO [67]. Except raw attention and Rollout, which are defined to be text-agnostic, the others are all text-specific, so we test the same image with two different text inputs on MS COCO. Our Grad-ECLIP demonstrates a strong ability of generating clear

and distinct text-specific heat maps, and gives reasonable explanation of verbs for interpreting CLIP. For example, the highlights for "holding" focus around the person's hands (the 5th row of Fig. 3i), while "standing" highlights the person's legs (the 6th row of Fig. 3i). We also notice that the sticks in the background are highlighted, which is probably because the sticks are regarded as "standing" in the snow.

In contrast to our methods, Grad-CAM and MaskCLIP can produce highlights on the explained object, but both also generate significant noise. CLIPSurgery tends to put high and coarse attention on the object region, but also contains background noises. M2IB and MaskCLIP fail when the texts are verbs ("holding" and "standing"), while RISE performs the worst with interpreting CLIP model. The results of GAME and Rollout, which are both based on self-attentions of the model, generate confusing heat maps due to the sparse attention between tokens in some layers.

### 4.1.2 Explanations on image-sentence pairs

The explanation map from Grad-ECLIP can also be generated from text encoder viewpoint. Using the gradient of matching score and the feature embeddings of word tokens, Grad-ECLIP can show the importance of each word in the given sentence

when matching with an image. Fig. 4 shows example visual and textual explanations for image-text pairs from MS COCO. Although Rollout and GAME can highlight important words in the sentence, Grad-ECLIP is the only one showing good correspondence between image attention regions and important words. From the explanation of the sentence, we can identify which words are more important for CLIP when matching with the specific image, and correspndingly the text-specific important regions on the image are shown in the visual explanation. This word importance visualization of the input text can be helpful when designing text prompts for image-text dual-encoders in practical applications.

### 4.1.3 Visualization examples on diverse image domains

We show the visualization comparison of different methods on the samples from different image domains, including the original ImageNet and ImageNet in different domains: rendition (ImageNet-R [68]), pencil sketch (ImageNet-Sketch [69]), natural adversarial example (ImageNet-A [70]), web images with captions (Conceptual Captions (CC) [71]), and chest x-ray with text (MSCXR [72]) in Fig. 5. For the image-caption pairs from the web-collected CC and chest X-ray data MSCXR, we generate explanations for both image and text encoder, and compare with the other methods that also provide text encoder explanations, including the raw attention, Rollout, GAME, M2IB.

Our Grad-ECLIP explanations provide interesting insights into how CLIP handles different image domains. In Fig. 5 (top), given a normal banana image and text "banana", Grad-ECLIP reveals that the yellow color is dominant to CLIP. However, when given a pencil sketch without color (ImageNet-Sketch), Grad-ECLIP reveals that CLIP looks at the curvature of the banana. For the color sketch of the banana (ImageNet-R), Grad-ECLIP shows that the color of the banana is mainly used, and not the black curved lines. Thus, from these examples, we may infer that CLIP prefers using the yellow color over the curved shape for matching with the "banana" text. With the sample of a web image and caption in the CC dataset, Grad-ECLIP generates clear and reasonable visual and the corresponding textual explanation map, showing that "cowboy", "horses" and "desert" are the main used concepts in the matching (from high to low importance).

Grad-ECLIP also provides interesting insights on how the original CLIP fails on novel domains. The last 2 rows of Fig. 5 show the explanations for chest x-ray images and text for the OpenAI CLIP model and a fine-tuned CLIP model (on MSCXR). The Grad-ECLIP explanation shows that the original CLIP uses the whole lobe to match with the words "defined" and "lobe". In contrast, the fine-tuned CLIP locates the actual anomaly and matches it with the text "defined opacities largely". The reason is that the fine-tuned model is trained to the specific domain that matches the X-ray and the illness location descriptions, while the original OpenAI CLIP model is more general and apparently "lobe" is the key word and the main object in the image-text pair.

### 4.1.4 Explaining ViT classifier and BLIP with Grad-ECLIP

Since our explanation method is designed for CLIP encoders, which are Transformer-based, our method can be easily adapted to generate visual explanations for other Transformer-based models. Here we adapt Grad-ECLIP to ViT-based classifier [73] and BLIP [8] to show the generality of our method.

For explaining ViT-based classifier, the classification score on the corresponding category is used to calculate the gradient and generate the heat map. From the examples for ViT classifier in
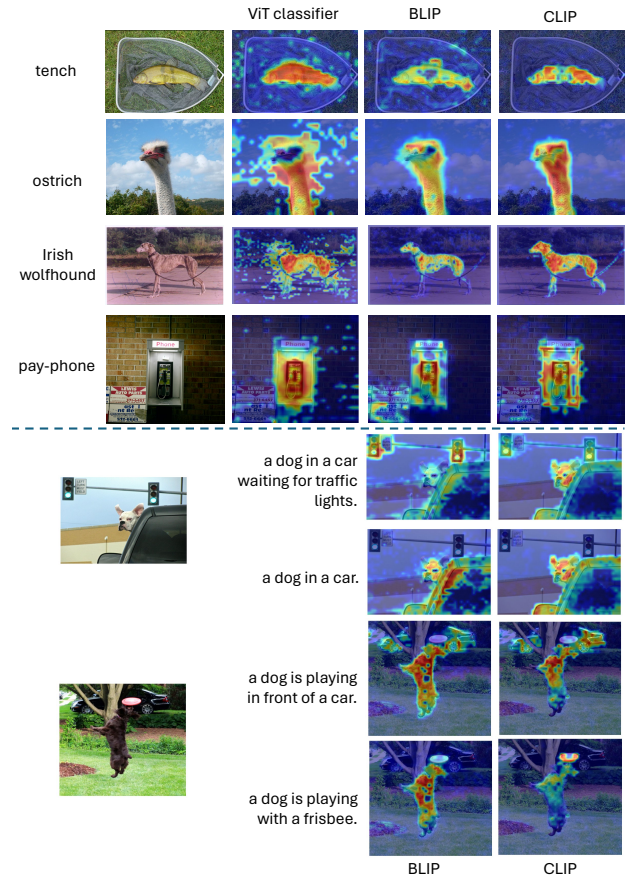


Fig. 6: Grad-ECLIP visual explanations of the ViT classifier and BLIP.

Fig. 6, we can see that although there are some noises on the background, Grad-ECLIP can well mark out the important region on the image for the specific class. As for the VLMs, shown by the visual explanation results presented in Fig. 6, when matching the same image-text pair, different models put attention on different regions. For example, BLIP notes the fins, while CLIP notes the fish body to match the image to "tench". When matching with the sentence "a dog is playing with a frisbee", BLIP puts more importance on the dog on the image, while CLIP shows places more importance on the frisbee.

Other VLMs like ALBEF [74] add additional attention layers after the encoders to fuse the image and text features, and thus our current method is not directly applicable since our method assumes that the last layer attention output has linear relationship with the final feature embedding. Our future work will investigate adapting our method to these modified ViT frameworks, e.g., the ALBEF model that uses cross attention to fuse image and text. Nonetheless, our ability to explain CLIP and other VLMs with similar architecture is significant considering that CLIP is by far the most widely used VLM.

### 4.1.5 Applying Grad-ECLIP to explain CNN-based CLIP

Although the methodology in §3.2 for Grad-ECLIP is derived from CLIP with Transformer architecture, here we show that our method is also applicable to CNN-based CLIP by using the attention layer in the final attention pooling. Figure 7 shows the visual explanation results for ResNet50-CLIP using Grad-ECLIP, and compared to other explanation methods that are compatible with CNN-based CLIP, including Grad-CAM, CLIPSurgery, and
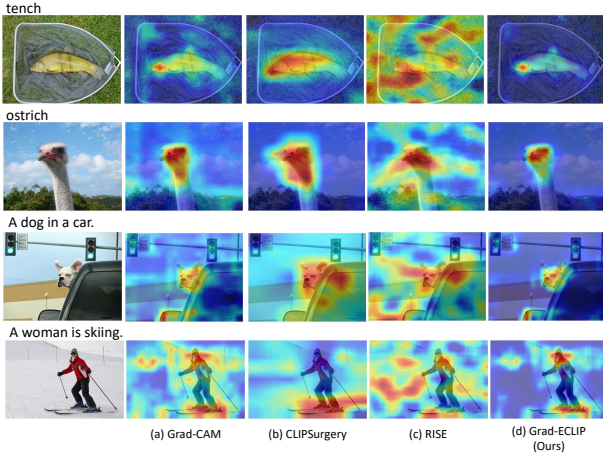
Fig. 7: Visual explanations of the CNN-based CLIP with ResNet50 architecture using Grad-ECLIP and other visual explanation methods.

RISE. From the examples, our Grad-ECLIP is able to generate clear and reasonable text-specific heat maps for interpreting CNN-based CLIP. In contrast to our method, Grad-CAM can produce similar highlight on the explained objects, but its heat map is more noisy in the background. Similar to its performance on ViT-based CLIP, CLIPSurgery generates rougher heat maps, which tend to put high values on a coarse region of the object. Meanwhile, RISE fails to explain CLIP model.

## 4.2 Quantitative evaluation

In this section, we perform quantitative evaluations of Grad-ECLIP comparing with baselines. In §4.2.1, the explanation faithfulness is evaluated by the Deletion and Insertion metrics [75, 34, 36, 37, 43], which are also called perturbation tests [15, 16]. Moreover, in §4.2.2, we evaluate localization ability, when considering each visualization as a soft-segmentation of the image, using PointGame [76, 77] and segmentation tests [15]. Finally, in §4.2.3, we evaluate the processing time of Grad-ECLIP compared with other visual explanation methods. In this section, in order to understand the effect of the spatial importance term in (19), we also present Grad-ECLIP without using the spatial weight $\lambda_i$, i.e., setting $\lambda_i = 1$, which is denoted as "w/o $\lambda_i$".

### 4.2.1 Deletion and Insertion

A faithful explanation method should produce heat maps highlighting the important content in the image that has greatest impact on the model prediction. Deletion (negative perturbation) replaces input image pixels by random values step-by-step with the important pixels removed first based on the ordering of the heat map values, while recording the drop in prediction performance. Insertion adds image pixels to an empty image step-by-step based on the heat map importance, and records the performance increase. For deletion, larger drops are better, while for insertion larger increase is better. We consider each step as 0.5% of number of image pixels, and record results for 100 steps. The model performance is measured using top-1 or top-5 zero-shot classification accuracy on the validation set of ImageNet [66] (ILSVRC) 2012, consisting of 50K images from 1000 classes. In particular, the perturbed image and each of the class names is fed into CLIP, and then classes with the highest scores are selected.

The insertion/deletion curves for top-1 accuracy are presented in Fig. 8, and the corresponding area under the curve (AUC) with top-1 and top-5 accuracy are presented in Tab. 1. Steeper drop of

performance with deletion steps corresponds to a lower deletion AUC, while quicker increase of performance with insert steps outputs a higher insertion AUC. Our method obtains the fastest performance drop for Deletion and largest performance increase for Insertion compared with most related works, showing that regions highlighted in our heat maps better represent explanations of CLIP. CLIPSurgery has comparable results to ours for Insertion, while performs poorly when evaluated with Deletion. The reason is that CLIPSurgery exhibits heat maps with nearly the same high values on the explained target region, so that the deletion operation fails to delete the most important pixels on the image at the beginning steps, which causes the deletion curve to decrease gradually, producing the high deletion AUC. Since CLIPSurgery can locate the explanation target with high values on the heat map, it performs well in the Insertion test. Our method without using the spatial importance (w/o $\lambda_i$) has slightly worse performance, but is still better than other baselines. As with [15, 16], we also use both the ground-truth class and the predicted class as the text prompt to generate heat maps, and our method is consistent with them, showing gains when using ground-truth text prompts.

We next evaluate the Deletion and Insertion performance for image and text retrieval tasks on the Karpathy's validation split of MS COCO. To evaluate the image explanations, image pixels are removed or added step-by-step based on the text-specific explanation heat map. With the modified image replacing the original image, we record the image and text retrieval performance (recall @ top-1 and top-5 matching) to draw the deletion and insertion curve. Then, the AUC results of Deletion and Insertion with text-specific image explanations are reported in Tab. 2. Grad-ECLIP surpasses the other methods on most metrics, which further demonstrates that our method produce high-quality visual explanation, regardless if the text is the class categories as in ImageNet or long captions as in MS COCO.

Finally, we evaluate *the faithfulness of our text explanations* using the *text version* of Deletion and Insertion metric, where words are deleted or inserted based on the order of importance in the text heat map. Using images and caption annotations in MS COCO Karpathy's split, we record the image-text retrieval performance for the modified caption, changing with total 5 steps with one word deleted/inserted at a time. The results in Tab. 3 show that Grad-ECLIP has the highest faithfulness (best deletion and insertion scores) compared with the other Transformer explanation methods. This demonstrates that Grad-ECLIP also has the excellent ability for image-specific text explanations.

### 4.2.2 Point Game and Segmentation Test

We next evaluate the localization ability of the visual explanations. We adopt the ImageNet-Segmentation (ImageNet-S) [78] validation set, which has segment annotations on 12,419 images of 919 categories from ImageNet. Point Game (PG) is a commonly used metric to evaluate the localization correctness of a visual explanation. PG counts a hit score if the location with the largest value in the visual explanation heat map lies within the object region, which can be defined by the class segmentation mask. Then the PG accuracy is measured by averaging over all samples. Since PG only considers the maximum point, but not the spread of the heat map, we also conduct energy-PG [36], which calculates the proportion of heat map energy within the ground-truth mask versus the whole map. Similar to the evaluation by [15, 16], regarding the heat maps as soft-segmentation results, we adopt

TABLE 1: Faithfulness evaluation of **image** explanation on the *ImageNet* validation set: AUC for Deletion and Insertion curves, based on Top-1 (@1) or Top-5 (@5) classification accuracy. Either the ground-truth or the prediction are used as the text input into CLIP. The second best is shown with underline.

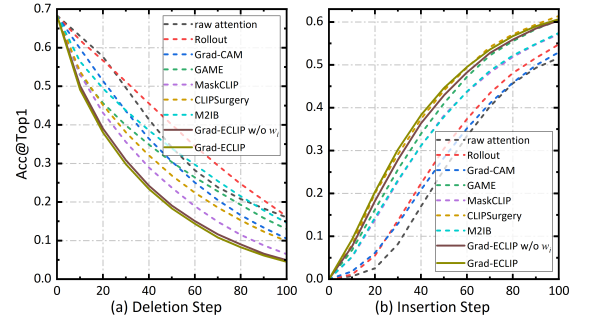| | Deletion↓ | | | | Insertion↑ | | | |
| | Ground-truth | | Prediction | | Ground-truth | | Prediction | |
| Method | @1 | @5 | @1 | @5 | @1 | @5 | @1 | @5 |
|---|---|---|---|---|---|---|---|---|
| raw attention | 0.3831 | 0.6239 | - | - | 0.2492 | 0.4195 | - | - |
| Rollout | 0.4082 | 0.6556 | - | - | 0.2803 | 0.4665 | - | - |
| Grad-CAM | 0.3417 | 0.5628 | 0.3518 | 0.5817 | 0.2682 | 0.4454 | 0.2526 | 0.4206 |
| GAME | 0.3356 | 0.5734 | 0.3497 | 0.5938 | 0.3611 | 0.5636 | 0.3425 | 0.5384 |
| MaskCLIP | 0.2848 | 0.4885 | 0.2886 | 0.4957 | 0.3335 | 0.5351 | 0.3275 | 0.5267 |
| CLIPSurgery | 0.3115 | 0.5235 | 0.3217 | 0.5412 | 0.3832 | **0.6021** | **0.3727** | 0.5719 |
| M2IB | 0.3630 | 0.5953 | 0.3633 | 0.5951 | 0.3351 | 0.5411 | 0.3347 | 0.5410 |
| Ours w/o $\lambda_i$ | <u>0.2535</u> | <u>0.4379</u> | <u>0.2634</u> | <u>0.4568</u> | <u>0.3715</u> | 0.5831 | <u>0.3528</u> | <u>0.5556</u> |
| Ours | **0.2464** | **0.4272** | **0.2543** | **0.4420** | **0.3838** | <u>0.5993</u> | <u>0.3672</u> | **0.5749** |



Fig. 8: Classification accuracy at Top-1 vs. (a) Deletion steps and (b) Insertion steps, on the *ImageNet* validation dataset with visual explanation heat maps from our Grad-ECLIP (solid) and other methods (dash).

TABLE 2: Evaluation of **image** explanation faithfulness on *MS COCO image-text retrieval (Karpathy's split)* validation dataset: AUC for Deletion and Insertion curves for performance on image retrieval (IR) and text retrieval (TR) tasks.

| | Deletion↓ | | | | Insertion↑ | | | |
| | IR | | TR | | IR | | TR | |
| Method | @1 | @5 | @1 | @5 | @1 | @5 | @1 | @5 |
|---|---|---|---|---|---|---|---|---|
| raw attention | 0.1708 | 0.3554 | 0.1923 | 0.3720 | 0.1247 | 0.2552 | 0.1544 | 0.2969 |
| Rollout | 0.1948 | 0.3946 | 0.2268 | 0.4238 | 0.1294 | 0.2932 | 0.1753 | 0.3503 |
| Grad-CAM | 0.1717 | 0.3502 | 0.2161 | 0.4008 | 0.1027 | 0.2216 | 0.1152 | 0.2327 |
| GAME | 0.1706 | 0.3552 | 0.1982 | 0.3800 | 0.1537 | 0.3083 | **0.2097** | 0.3735 |
| MaskCLIP | 0.1321 | 0.2841 | **0.1516** | 0.2949 | 0.1423 | 0.2953 | 0.1891 | 0.3514 |
| CLIPSurgery | 0.1794 | 0.3652 | 0.2381 | 0.4292 | 0.1419 | 0.2941 | 0.1771 | 0.3384 |
| M2IB | 0.1797 | 0.3671 | 0.2057 | 0.3905 | 0.1469 | 0.3004 | 0.2058 | 0.3691 |
| Ours w/o $\lambda_i$ | 0.1390 | 0.2940 | 0.1827 | 0.3386 | 0.1403 | 0.2895 | 0.1735 | 0.3279 |
| Ours | **0.1246** | **0.2670** | 0.1550 | **0.2933** | **0.1576** | **0.3203** | 0.2056 | **0.3761** |

TABLE 3: Evaluation of **text** explanation faithfulness on *MS COCO image-text retrieval (Karpathy's split)* validation dataset: AUC for Deletion and Insertion curves with reporting image retrieval (IR) and text retrieval (TR) performance.

| | Deletion↓ | | Insertion↑ | |
| Method | IR | TR | IR | TR |
|---|---|---|---|---|
| raw attention | 0.2843 | 0.4917 | 0.0065 | 0.0328 |
| Rollout | 0.1221 | 0.2389 | 0.1052 | 0.2070 |
| GAME | 0.1083 | 0.2084 | 0.1146 | 0.2301 |
| M2IB | 0.2139 | 0.4256 | 0.0063 | 0.0375 |
| Ours w/o $\lambda_i$ | 0.1116 | 0.2113 | 0.1123 | 0.2361 |
| Ours | **0.0996** | **0.1770** | **0.1292** | **0.2536** |

pixel accuracy (Pixel Acc.), average precision (AP), and averaged mask intersection over union (maskIoU) as additional metrics.

The results for localization are shown in Tab. 4. Both versions of Grad-ECLIP significantly outperform other explanation methods on PG and energy-PG, which demonstrates that Grad-ECLIP can well reveal that the important pixels for CLIP are inside the object region. Comparing Grad-ECLIP with and without $\lambda_i$, Grad-ECLIP without $\lambda_i$ obtains relatively higher performance on pixel accuracy and maskIoU, since heat maps that contain more high-value pixels within the ground-truth mask have advantage on these two metrics. In Fig. 9(b,c), using $\lambda_i$ reduces the values on the mask while removing the surrounding noise. Due to the similar reason, CLIPSurgery obtains higher pixel accuracy and maskIoU, since it tends to put high heat map values on all the pixels of object region, and gets higher score when aggregating the heatmaps inside the object mask in these two evaluations. However, the lower PG, energy-PG and AP demonstrate that there are more high values generated outside of the object boundary. Better segmentation does not necessarily result in faithful explanations, in terms of both insertion and deletion metrics, as indicated in Table 1.

TABLE 4: Evaluation of localization ability using the Point Game (PG and energy-PG) and Segmentation test (Pixel Acc., AP and MaskIoU) on the *ImageNet-S* validation dataset.

| Method | PG | energy-PG | Pixel Acc. | AP | maskIoU |
|---|---|---|---|---|---|
| raw attention | 0.1219 | 0.1321 | 0.0278 | 0.2877 | 0.0013 |
| Rollout | 0.1375 | 0.2835 | 0.2524 | 0.3345 | 0.011 |
| Grad-CAM | 0.1845 | 0.3154 | 0.5457 | 0.4050 | 0.1251 |
| GAME | 0.4706 | 0.4438 | 0.4765 | 0.4072 | 0.089 |
| MaskCLIP | 0.4041 | 0.1408 | 0.718 | 0.4557 | 0.2481 |
| CLIPSurgery | 0.5759 | 0.3983 | **0.7546** | 0.4608 | **0.3471** |
| M2IB | 0.264 | 0.3557 | 0.6194 | 0.4003 | 0.1474 |
| Ours w/o $\lambda_i$ | <u>0.8356</u> | <u>0.4409</u> | <u>0.7365</u> | <u>0.5163</u> | <u>0.3314</u> |
| Ours | **0.8899** | **0.5997** | 0.7056 | **0.5662** | 0.2869 |

operations. Note that for gradient-based methods, the backpropagation does not need to go all the way to the input layer, but stops at an intermediate upper layer, and thus the extra computation required is not much. RISE needs the longest processing time, which is a common drawback of perturbation-based methods.

### 4.3 Ablation study

In this section, we conduct ablation studies to illustrate the influence of the proposed loosened spatial weight (§4.3.1), the number of involved layers (§4.3.2) and multi attention heads (§4.3.3) in the calculation of Grad-ECLIP.

#### 4.2.3 Processing time comparison

In Tab. 5, we show the average processing time per image, which counts the total duration from inputting the image and text into CLIP to obtaining the explanation map. Since the gradient can be easily and quickly obtained through the autograd function of Pytorch, both our method and Grad-CAM takes similar processing time as the raw attention and MaskCLIP, which obtain their heat maps from the forward pass of the model and some other minor

#### 4.3.1 Effect of the loosened spatial weight

Here we conduct ablation study on the $\lambda_i$ in Eq. 20 to show the effect of the proposed spatial weight. We consider two versions

TABLE 5: Comparison of the average processing time (on RTX3090 GPU) per image for generating the explanation map.

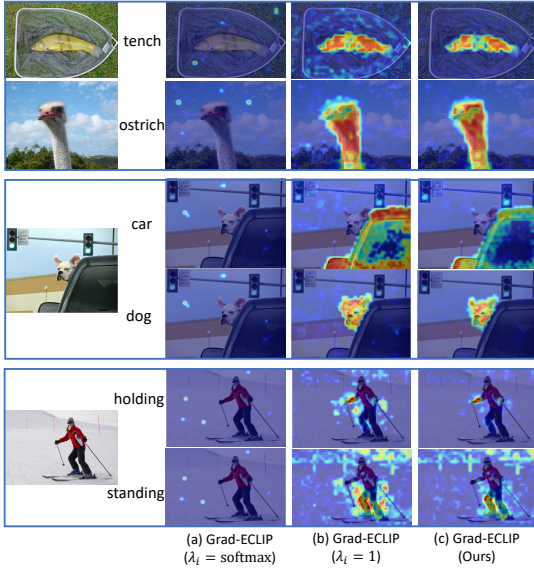| Method | raw attention | Rollout | Grad-CAM | GAME | MaskCLIP | CLIPSurgery | M2IB | RISE | Grad-ECLIP(Ours) |
|---|---|---|---|---|---|---|---|---|---|
| time (s/img) | 0.0117 | 0.0298 | 0.0114 | 0.0228 | 0.0117 | 2.9423 | 0.5781 | 6.2376 | 0.0165 |



Fig. 9: The effect of spatial weight $\lambda_i$. We compare the explanation maps of Grad-ECLIP with a version without the proposed $\lambda_i$ (denoted as "w/o $\lambda_i$", which replaces the proposed spatial weights with (a) $\lambda_i = \text{softmax}$, (b) $\lambda_i = 1$), and (c) the full Grad-ECLIP.



Fig. 10: The **image** visual explanations generated when aggregating over $N$ layers of the image transformer encoder.



Fig. 11: The **textual** explanations generated when aggregating over $N$ layers of the text transformer encoder.

of Grad-ECLIP with modified spatial importance $\lambda_i$: the first version uses the softmax attention rather than 0-1 normalization (denoted as $\lambda_i = \text{softmax}$); the second version removes the spatial importance completely by setting $\lambda_i = 1$.

A comparison of visualizations is presented in Fig. 9. Compared with the heat maps generated by the full Grad-ECLIP in Fig. 9c, the version that removes spatial importance altogether ($\lambda_i = 1$) contains more noise near object boundaries and on the background (Fig. 9b), but are otherwise consistent with full Grad-ECLIP. The result of using $\lambda_i = \text{softmax}$ (Fig. 9a) is equivalent to raw attention (Fig. 3a) due to the output of the $\text{softmax}$ being extremely sparse. We also provide the quantitative comparisons of Grad-ECLIP using $\lambda_i = 1$ in the quantitative evaluation oin §4.2.1 and §4.2.2, denoted as "w/o $\lambda_i$".

### 4.3.2 Effect of number of layers on visual explanation

As introduced in §3.2, the explanation can be aggregated over all the layers in Transformer by recursively processing each layer with Eq. 19. In this section, we conduct the experiments to discuss the influence of using different number of layers to generate heat maps for image and text.

With different layer number $N$, the visualizations with specific texts are shown in Fig. 10, and the corresponding caption explanations are shown in Fig. 11. $N = 1$ means the visualization is generated only with the final Transformer layer, while $N = 12$ means all the layers are involved. The image explanations become worse when increasing the number of layers involved, since the features in lower layer may introduce more noise to the heat map. Therefore, it is the best to just use the last layer in the calculation of image visual explanation, where this conclusion is consistent with the classical gradient-based CAM methods.

As for the text explanation, there is no obvious difference of the visualization quality, since the highlights are basically focusing
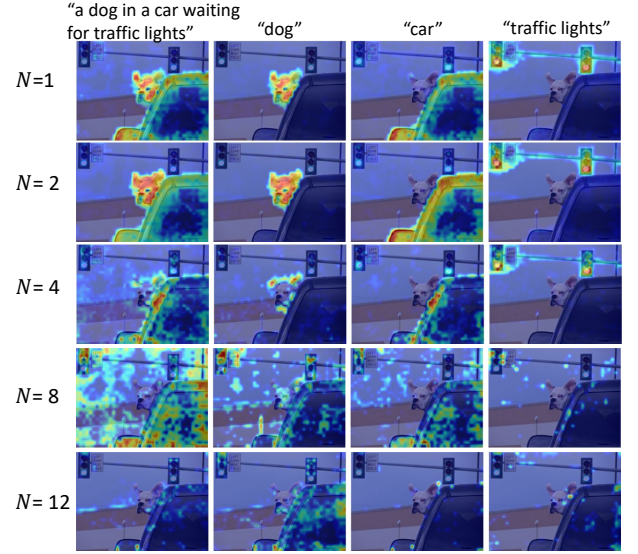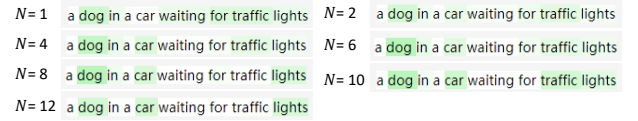
on "dog", "car" and "traffic lights" with some minor variations. Therefore, we perform the Deletion and Insertion experiments as in §4.2 on the text explanation maps based on different number of layers $N$. The results are shown in the following Table 6. The explanation faithfulness has the trend that it first increases with more layers used and then goes down with the lower-layer features involved ($N > 8$). Therefore, we aggregate the last eight layers maps for interpreting the text encoder in our experiments.

### 4.3.3 Effect of multi attention heads on visual explanation

As mentioned in (7) of §3.2, for producing the Grad-ECLIP visual explanation, we set CLIP to perform the forward pass with a single head in the attention layer instead of the original multi-head attention layer. In Fig. 12, we show the visualization of explanation maps when using multi-head attention layers, compared to using

TABLE 6: The **text** explanation faithfulness vs. the number of transformer layers aggregated for the explanation. Evaluating on *MS COCO image-text retrieval (Karpathy's split)* validation dataset: AUC for Deletion and Insertion curves with reporting image retrieval (IR) and text retrieval (TR) performance.

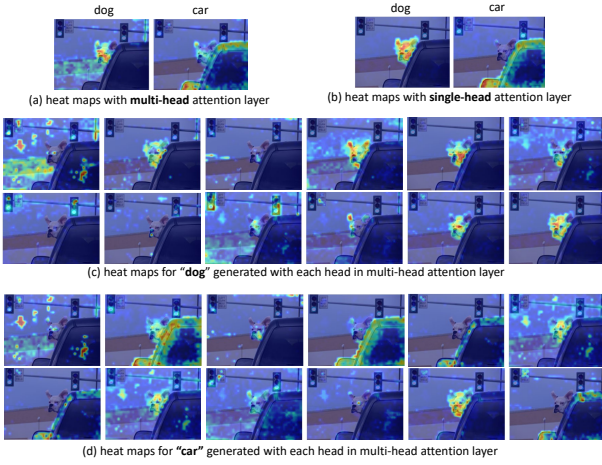| $N$ | Deletion↓ | | Insertion↑ | |
|---|---|---|---|---|
| | IR | TR | IR | TR |
| 1 | 0.1118 | 0.2087 | 0.1059 | 0.2196 |
| 2 | 0.1021 | 0.1826 | 0.1186 | 0.2351 |
| 4 | **0.0995** | 0.1786 | 0.1242 | 0.2428 |
| 6 | 0.0989 | **0.1761** | 0.1273 | 0.2490 |
| 8 | 0.0996 | 0.1770 | **0.1292** | **0.2536** |
| 10 | 0.1008 | 0.1843 | 0.1288 | 0.2472 |
| 12 | 0.1095 | 0.2087 | 0.1219 | 0.2364 |

Fig. 12: The visual explanation maps with using (a) multi-head attention layer; (b) single-head attention layer; (c) each head in multi-head attention layer for text "dog"; (d) each head in multi-head attention layer for text "car"
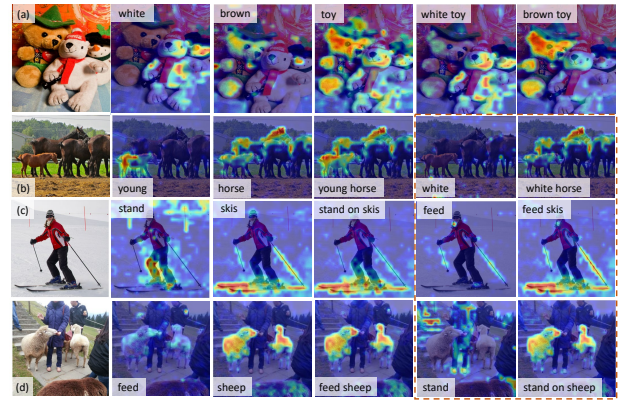


Fig. 13: Visual explanation heat maps generated for single words and word phrases using Grad-ECLIP on CLIP. The dashed box contains examples where the text does not match the image.

a single head. Comparing Fig. 12 (a) and (b), using multi-head attention results in some surrounding context information is also highlighted with the explained object.

We further produce the heat maps for *each* attention head, using the $q \in \mathbb{R}^{D/12}$, $k \in \mathbb{R}^{D/12}$, $v \in \mathbb{R}^{D/12}$, and attention output $o_{cls} \in \mathbb{R}^{D/12}$, where $D$ is channel number before going into multi heads, and visualize them in Fig. 12 (c) for the target "dog", and (d) for the "car". The visual explanation in each head highlights different regions, not only seeing the target object. We can infer that the channels assigned to each heads can preserve different information, and the softmax inside each head helps the model to encode more context information. In contrast, with the single head setting, the softmax is performed over all channels, which selects out the most important information, and our explanation method can show the model's attention on the specific explained target, as shown in Fig. 12 (b).

# 5 ANALYSIS OF CLIP USING GRAD-ECLIP

Useful explanation methods can be used to identify failure modes, establish appropriate users' confidence and give insight to developers to improve models. Therefore, in this section, we use the visual explanation maps and textual explanations generated by Grad-ECLIP to give examples of how to explore the mechanism in text and image matching, and analyze the strengths, weaknesses, and preference of CLIP model. We hope that our explanation tool can help researchers discover more interesting properties of pretrained image-language models, and inspire further development of these models. Here we analyze three aspects of CLIP: 1) the decomposition and addibility in image-text matching (§5.1); 2) types of attributes that can be identified by CLIP (§5.2); 3) the concreteness/abstractness of words learned by CLIP (§5.3).

## 5.1 Concept decomposition and addibility in image-text matching

Examining the visualizations shown in Fig. 3h, CLIP can well recognize the single concepts (nouns) and has good attention about actions (verbs). An interesting question is how does it process the combination of words, *e.g.*, adjective and noun, verb and noun? To examine the working function of phrase matching, we conducted

experiments comparing the explanation heat maps for single words and combined phrases using Grad-ECLIP.

The results are shown in Fig. 13. Considering adjective-noun combinations in (a), the highlights are put on all three toys when matching with "toy", and CLIP can successfully highlight the correct toy when the color adjective is included in the text. In the case of "young horse" in (b), the other horses are still highlighted, while the highlights on the young one is strengthened by adding the attribute "young". The examples of verb-noun cases in (c) and (d) also show similar addibility pattern on the heat maps: (c) with the verb "stand", the region of person's leg is highlighted along with the "skis"; (d) with the verb "feed", the people's hands are also highlighted together with sheep. We also show some non-existent concepts or strange word combinations in the dashed box of Fig. 13, e.g., "white horse" in (b), "feed skis" in (c), "stand on sheep" in (d). In these cases, the visualization shows that CLIP will mainly focus on the reasonable part of the concept, such as "horse", "skis" and "sheep". For the non-existent "white" concept in image (b), the visual explanation does not highlight anything.

Therefore, we infer that when processing the matching of image and phrases, the model has the ability of decomposition and addibility of different concepts. This can help the model to generalize to different scenarios and could be the source of the strong zero-shot ability of CLIP.

## 5.2 Diagnostics on attribution identification

In Fig. 13(a), we see that CLIP has an ability to distinguish color attributes, and mark out the corresponding regions on image. To explore further, we conduct an experiment to test CLIP's ability to identify different types of object attributes. We adopt an example image from CLEVR [79], a diagnostic dataset for visual reasoning, and visualize image-text matching with various attributes: shape (sphere, cylinder, cube), material (metal, matte, plastic), color (red, yellow, blue), size (big, small), position (left, right).

Fig. 14 shows the visual explanation heat maps generated with each image-attribution pair. We have the following findings: 1) for shape and material, the heat maps can show partial correct attention with some obvious objects, such as the metal sphere for "sphere" and the highlighted cylinder and cube for "matte". However, there are also false positive and false negative errors in (a) and (b). Thus, CLIP possess a certain but limited knowledge about object shapes and materials. 2) For the color attribute in row (c), the results further verify that the model can have good ability to distinguish different colors. 3) For comparative attributes, size
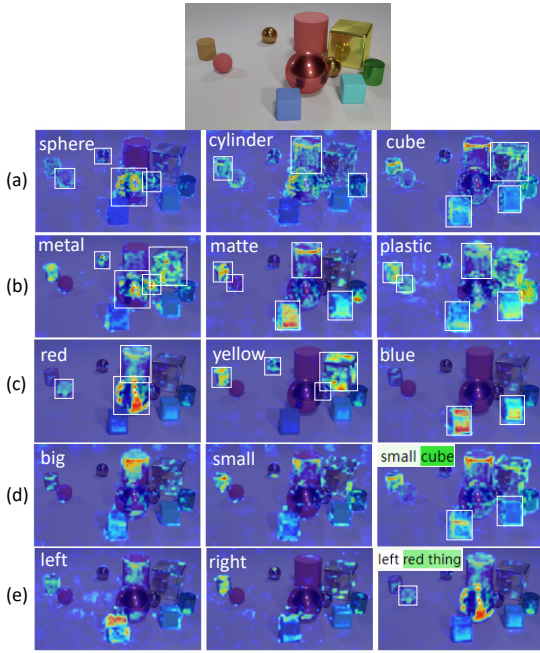
Fig. 14: Visual explanations on image matching with different kinds of attributions: (a) shape; (b) material; (c) color; (d) size; (e) position. For visualization, the ground-truth corresponding to the text prompt are outlined with white boxes, except for cases involving relative adjectives, *e.g.* "big", "small", "left", "right". The text explanation maps are also shown for "small cube" and "left red thing" combinations.
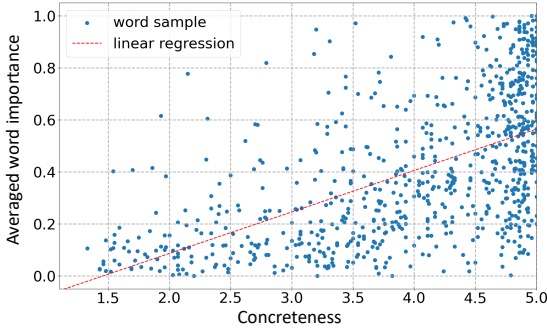


Fig. 15: The average word importance (via Grad-ECLIP) vs. word concreteness for the top-1000 most frequent words in the MS COCO Karpathy's validation split. The red dash line shows the linear regression result over the word samples; $r^2 = 0.32, p < 0.001$.

(big or small) in (d) and position (left or right) in (e), the visual explanations also show that CLIP produces some erroneous results. For example, there are little differences between the heat maps of "small cube" and "cube" in (a), or "left red thing" and "red" in (c), which demonstrates that the word "cube" and "red" take the major role in the matching. This is also confirmed by the text heat maps in the figure.

Overall, from the above analysis, we infer that CLIP has advantages with common perceptual attributes like color, but cannot well handle physical attributes like shape and material, and is weak at grounding objects with comparative attributes, like size and position relationships. Related to the addibility of concepts in the §5.1, it is reasonable to expect that attributes that have concrete visual appearance, e.g., color, will contribute more to the matching score, compared with the abstract comparative attributes.

## 5.3 Relationship between word importance and concreteness

From §5.2, we have inferred that the concrete visual concepts (e.g., "red" and "blue") contribute more to the image and text matching than the abstract attributes (e.g., "left" and "right"). Therefore, based on the text explanation obtained from Grad-ECLIP, we further explore the relationship between word importance in matching and word concreteness, and analyze which type of concepts and words (concrete vs. abstract) that the CLIP model has actually learned and uses most often for matching. For an image-text pair, we calculate the textual explanation from Grad-ECLIP, and then obtain the importance value of each word by normalizing such that the maximum word importance value is 1 in the sentence. We then calculate the average word importance on the top-1000 most frequent words in the MS COCO caption (Karpathy's split) validation set. For the word concreteness, we adopt the open sourced database from [80], which provides the concreteness for 40,000 common English words measured by human rating. The concreteness is a value from 1 to 5, where 5 means the most concrete and 1 means the most abstract.

Based on the 1000 selected words, Fig. 15 presents a scatter plot of average word importance versus concreteness value. We perform linear regression analysis on this data (red dashed line), and the regression result was statistically significant ($r^2 = 0.32$, $p < 0.001$). The scatter distribution and linear regression result reveal that CLIP places higher word importance to more concrete words, and vice versa, less word importance on more abstract words. *Therefore, the words that CLIP has learned for matching are biased towards concrete words.*

CLIP's learned bias towards concrete words could be due to frequency bias in the training set, i.e., concrete words could appear more frequently during training. To investigate this possibility, we attempt to count the number of occurrences of the selected words in CLIP's training set WebImageText [1]. However, since WebImageText is not publicly available, we instead compute these statistics from the OpenWebText [81] dataset, which follows the same methodology to reproduce the data characteristics and structure of the WebImageText corpus.

The relationship between average word importance and the word frequency in the training corpus is plotted in Fig. 16, with different sample colors representing the levels of concreteness. There is no obvious relationship between the word attention and the frequency. Many high frequency words obtain a low word importance when the concreteness value is very low, and on the opposite, low frequency words of concrete concepts may obtain a high word importance. Fig. 16 (right) shows the zoom-in on the samples in the blue box, which have high average word importance and in the purple box with low average importance. Comparing these two boxes, the words with concrete visual meaning, especially nouns, are indeed more likely to obtain higher importance than the abstract words. Therefore, the analysis using Grad-ECLIP text explanation further supports the conclusion from §5.2 – *concrete visual concepts are learned better and contribute more to the CLIP image-text matching score than abstract concepts, and this phenomenon is not due to word frequency bias.*

## 6 APPLICATION OF GRAD-ECLIP TO BOOST FINE-GRAINED UNDERSTANDING OF CLIP

CLIP has been shown to have limitations in understanding fine-grained details, such as poor region recognition when using its
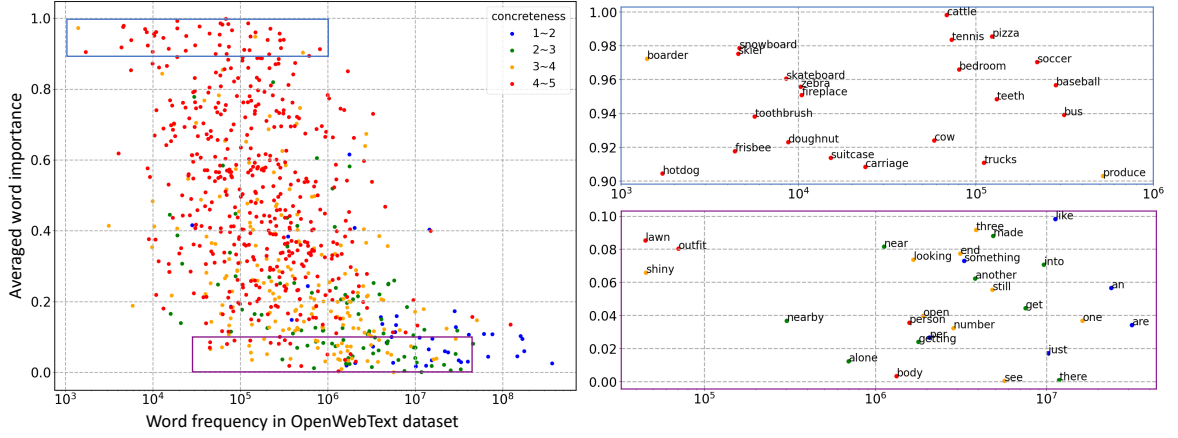
Fig. 16: The average word importance (via Grad-ECLIP) vs. the word frequency in the OpenWordText dataset for the top-1000 most frequent words in the MS COCO Karpathy's validation split. The colors represent the concreteness level of each word according to [80]. (right) zoom-in of the blue and purple boxes to show the word examples with high word importance (blue box) and low word importance (purple box).
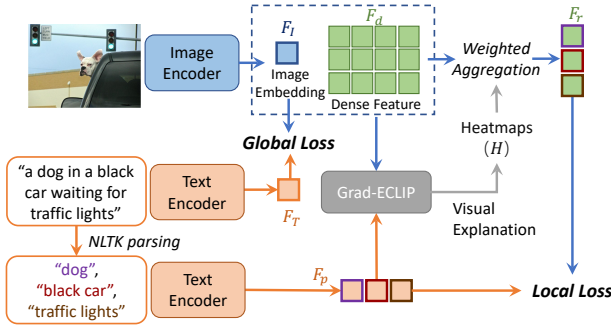


Fig. 17: Overview of the proposed fine-grained fine-tuning of CLIP using Grad-ECLIP. Multiple phrases or words representing objects (e.g., "dog", "black car" and "traffic lights") are separated out by parsing the input caption via the Natural Language Toolkit (NLTK). With the visual explanation heat maps generated by Grad-ECLIP, object-specific region feature embeddings ($F_r$) are obtained through weighted aggregation of the image dense feature ($F_d$). Finally, the CLIP fine-tuning is composed of two losses: the global loss as in the original pre-training of CLIP, and the local loss, which aligns the image region features ($F_r$) and the corresponding phrase features ($F_p$).

dense features, due to the pre-training focusing on matching the whole image (via the $[cls]$ token) to a text description. Our Grad-ECLIP generates text-specific explanation maps for the image encoder, which highlights the detailed regions on the input image when CLIP matches the image-text pair. Therefore, aiming to boost the fine-grained understanding ability of CLIP, here we propose a fine-tuning method for CLIP that adopts Grad-ECLIP to indicate the fine-grained alignment between image regions and corresponding textual attributes. As a result, the fine-tuned CLIP has a representation space where visual and semantic features are both globally and locally aligned, which significantly improves its region aware ability for dense prediction tasks, as well as maintaining its performance on image-level tasks.

## 6.1 Grad-ECLIP-based fine-grained fine-tuning of CLIP

The framework of our proposed method is shown in Fig. 17, where the fine-tuned model includes an image encoder $\mathcal{I}(\cdot)$ and a text encoder $\mathcal{T}(\cdot)$, which is the same model architecture as the original CLIP. In the next two subsections, we will introduce the global loss and local loss adopted in the fine-grained fine-tuning, respectively.

### 6.1.1 Global loss

The input of CLIP is a batch of image-text pairs $\{(I_b, T_b)\}_{b=1}^B$. After passing through the image encoder and text encoder, the model outputs the corresponding global image and text embeddings $\{(F_{I_b}, F_{T_b})\}_{b=1}^B$, respectively. The fine-tuning adopts the same global contrastive learning as in the pre-training, which realizes the instances-level alignment and helps to maintain the model's ability for multi-modal global representation. The cosine similarity $S(F_I, F_T)$ between image embedding $F_I$ and text embedding $F_T$ is calculated as in Eq. 1. The constrastive loss is applied to push CLIP to learn the global representations by maximizing the cosine similarities of the corresponding image and text embeddings, while minimizing the cosine similarities of other non-paired ones in the batch, which is defined as:

$$
L_{global} = -\frac{1}{2B} \sum_{b=1}^{B} \left( \log \frac{\exp\left(S(F_{I_b}, F_{T_b})/\tau\right)}{\sum_{b'=1}^{B} \exp\left(S(F_{I_b}, F_{T_{b'}})/\tau\right)} + \right.
$$
$$
\left. \log \frac{\exp\left(S(F_{T_b}, F_{I_b})/\tau\right)}{\sum_{b'=1}^{B} \exp\left(S(F_{T_b}, F_{I_{b'}})/\tau\right)} \right), \quad (21)
$$

where $\tau$ is the trainable temperature parameter.

### 6.1.2 Local loss

The local loss is based on the extracting the dense feature map from the image and performing fine-grained matching of region features to the corresponding text descriptions.

**Image dense feature.** Following [19, 28] we extract the dense feature map of the input image from a ViT-based encoder by to slightly modifying the last transformer layer to keep the projection and norm layers, and discard the self-attention. This modification is experimentally shown to be capable of preserving more spatial detailed features in the output token embeddings. Specifically, in the last transformer layer, with the input $x = (x_{cls}, x_1, ..., x_{h \times w})$ comprising a $[cls]$ embedding and $h \times w$ spatial token embeddings, the output of the attention layer is $o = v = \mathcal{LP}(x)$ instead of the $o = \mathcal{A}(x) = \text{softmax}(\frac{qk^\top}{\sqrt{C}})v$ in (7). Then, the $[cls]$ embedding is removed and the final spatial token embeddings are reshaped into an $h \times w$ image dense feature map $F_d$, from which we extract fine-grained representations for specific image regions.

**Image-text fine-grained matching.** Next we design a region feature and text matching scheme based on Grad-ECLIP to automatically obtain region-text alignment pairs. Note that our method

does not require any manual or network-based region proposal or label annotations. For the caption $T$ in each image-text pair $(I, T)$, we use the Natural Language Toolkit (NLTK) [82] to parse and extract the phrases that contain object concepts, by setting the separation and selection rules as "adjective + noun". For the example in Fig. 17, the NLTK extracts "dog", "black car", and "traffic lights" from the input text "a dog in a black car waiting for traffic lights". Then, these extracted words or phrases $\{p_t\}_{t=1}^n$, where $n$ is the maximum number of extracted concepts from each caption, are sent to the text encoder, resulting in a set of phrase embeddings $\{F_{p_t}\}_t$. A region embedding $F_{r_t}$ is then calculated for each phrase embedding. Specifically, the phrase embeddings $F_{p_t}$ is used to calculate the cosine similarity with the image embedding $F_I$. Grad-ECLIP is then applied to the calculated score to obtain a heat map $H_t$ according to the procedure in §3.2, which reveals the important spatial locations for matching with this specific phrase. Therefore, we adopt the explanation heat maps as weights for aggregating the image dense feature $F_d$, resulting in the region embedding for the phrase: $F_{r_t} = \sum_{hw} H_t \cdot F_d$, where $\sum_{hw}$ is the sum operator over spatial coordinates.

Finally, for the region embeddings and corresponding phrase embeddings, we adopt focal loss [64] to match the positive pairs and distance the negative pairs:

$$
\begin{aligned}
L_{local} = & -\sum_t \left(1 - S\left(F_{r_t}, F_{p_t}\right)\right)^2 \log S\left(F_{r_t}, F_{p_t}\right) \\
& -\sum_t \sum_{t' \neq t} S\left(F_{r_t}, F_{p_{t'}}\right)^2 \log \left(1 - S\left(F_{r_t}, F_{p_{t'}}\right)\right),
\end{aligned}
\tag{22}
$$

where $S$ represents the cosine similarity function, and $t$, $t'$ means the $t$-th and $t'$-th phrase in the same batch. By adding the local loss to the global loss as the total loss, our proposed fine-grained fine-tuning successfully boosts the representation alignment between image region and corresponding textual concepts, while maintains the image-level performance at the same time.

## 6.2 Experiments with fine-grained fine-tuning CLIP

In these experiments, we use our proposed Grad-ECLIP to enable fine-tuning of CLIP to enhance the fine-grained understanding.

### 6.2.1 Experiment settings

The experiments are conducted based on the pre-trained models from EVA-CLIP [83] considering its high efficiency and capacity, followed the previous work [28]. For the fine-tuning experiments, unless otherwise specified, we adopt the Conceptual Caption (CC3M) dataset [71] as the training set, which collects about 3 million image-text pairs from the internet. During fine-tuning, the input image size is 224x224, which is the same as the original pre-training of CLIP. Two RTX 6000 Ada are used with batch size 64 on each, learning rate of 1e-5, and weight decay of 0.1.

### 6.2.2 Experiment results

**Evaluating the fine-grained representation.** To evaluate the dense representation ability of the fine-tuned CLIP, we use the mean accuracy (mAcc) of classifying region boxes annotated in the val2017 split of MS COCO and panoptic masks (including "things" with 80 classes and "stuffs" with 91 classes) annotated in MS COCO Panoptic dataset [84]. For classification, RoI pooling is used to extract region box embeddings, while mask pooling is used to extract the mask embeddings from the image dense feature

maps. The classification is performed by selecting the highest score when matching with the text embeddings of the classes.

The results are shown in Tab. 7. Compared with the pre-trained CLIP base model, fine-tuning with just the global loss on CC3M dataset can slightly increase the zero-shot classification performance – this setting is equivalent to continuing to train CLIP in the original way, which we denote as *ordinary FT* (fine-tuning) in the following descriptions. When we use both the global and local loss with the help of our Grad-ECLIP explanation maps, the fine-tuned model's performances on region classification obtained significant improvements, for both boxes and masks. In Tab. 7, we also list the absolute increases of each metric compared with ordinary FT (global loss only) using the red values, which shows the effectiveness of the local loss on boosting the fine-grained representation of CLIP ViT model.

To demonstrate the effectiveness of our Grad-ECLIP, we also conduct the fine-tuning with the local loss using two other plug-in and low computation cost explanation methods, Grad-CAM and MaskCLIP. With the heat maps from these two methods, the region-aware matching also obtains obvious performance improvements compared with the ordinary FT, but there is still a gap compared with using our Grad-ECLIP. The comparison results further demonstrate that the high-quality and accuracy of the visual explanation maps generated by Grad-ECLIP.

**Application to open-vocabulary detection (OVD) task.** We adopt the fine-tuned CLIP as the backbone for OVD to verify the fine-grained understanding by a down-stream localization task. Following the previous work CLIPSelf [28], which is a state-of-the-art CLIP fine-tuning scheme for increasing the fine-grained region representation via self-distillation, we build open-vocabulary object detectors based on the F-ViT [28] architecture, which is a two-stage detector using a frozen CLIP ViT as the backbone. In the CLIP fine-tuning stage, we adopt the image-text pairs in CC3M or MS COCO Karpathy train set as the inputs, and for MS COCO Karpathy train set, we filter out the same image samples that are also in the OVD val set to ensure there is no risk of label leakage. The OVD models are trained on the OV-COCO benchmark [88], and we use AdamW optimizer with batch size of 64, learning rate of 1e-4, and weight decay of 0.1. For evaluation, we report box AP (average precision) at IoU (Intersection over Union) of base, novel and all categories as with previous works [13, 56, 57, 27, 28]. Since there is no extra region box annotations used in our fine-tuning, for fair comparison, we implement fine-tuning of the CLIPSelf version using image patch distillation, which also has no extra region proposal annotation requirements.

The results are presented in Tab. 8. F-ViT is the baseline that initializes the detector backbone with the original pre-trained CLIP model, and other versions with † initialize the backbone with the CLIP models fine-tuned by various methods and datasets. With the ViT-B/16 backbone, ordinary FT produces similar performance as the baseline, while our FT significantly improves the OVD results, especially on the novel categories. Since the base categories have explicit annotated bounding boxes and labels during OVD training, the performance on the unseen novel categories better illustrates the fine-grained understanding ability brought by the CLIP model. In contrast to CLIPSelf where the downstream OVD performance is much influenced by the dataset used for fine-tuning, our method obtains better performances regardless of the fine-tuning dataset (CC3M or MS COCO), which indicates that our method can effectively and stably boost the dense representation ability of

TABLE 7: Evaluating the fine-grained representation of the fine-tuned CLIP via zero-shot classification on the MS COCO validation dataset. We report the Top-1 and Top-5 mean accuracy on both object bounding boxes and panoptic masks (thing and stuff). Besides our Grad-ECLIP, we also adopt Grad-CAM and MaskCLIP to produce explanation map in the calculation of local loss for comparison. The red arrow shows the performance improvement brought by the local loss based on our Grad-ECLIP, compared with just using the global loss, which equivalent to ordinary fine-tuning. The gray row is the baseline CLIP before fine-tuning.

| | | Fine-tuning | | | Boxes | | Thing Masks | | Stuff Masks | |
| Method | Model | Global Loss | Local Loss | Explanation Map | Top1 | Top5 | Top1 | Top5 | Top1 | Top5 |
|---|---|---|---|---|---|---|---|---|---|---|
| CLIP | ViT-B/16 | - | - | - | 41.4 | 63.6 | 30.6 | 53.8 | 13.9 | 36.6 |
| Fine-tuned CLIP | ViT-B/16 | $\checkmark$ | - | - | 42.9 | 64.8 | 32.9 | 56.4 | 14.7 | 38.7 |
| | | $\checkmark$ | $\checkmark$ | Grad-CAM | 54.2 | 74.7 | 46.5 | 69.8 | 13.2 | 42.2 |
| | | $\checkmark$ | $\checkmark$ | MaskCLIP | 54.3 | 75.5 | 47.4 | 70.9 | 17.0 | 47.9 |
| | | $\checkmark$ | $\checkmark$ | Grad-ECLIP (Ours) | $57.3_{\uparrow 14.4}$ | $78.3_{\uparrow 13.5}$ | $49.3_{\uparrow 16.4}$ | $72.2_{\uparrow 15.8}$ | $18.3_{\uparrow 3.6}$ | $51.1_{\uparrow 12.4}$ |
| CLIP | ViT-L/14 | - | - | - | 58.1 | 78.9 | 49.8 | 72.6 | 13.1 | 33.9 |
| Fine-tuned CLIP | ViT-L/14 | $\checkmark$ | - | - | 62.6 | 83.1 | 54.7 | 77.5 | 16.2 | 39.2 |
| | | $\checkmark$ | $\checkmark$ | Grad-ECLIP (Ours) | $71.7_{\uparrow 9.1}$ | $89.7_{\uparrow 6.6}$ | $63.4_{\uparrow 8.7}$ | $85.6_{\uparrow 8.1}$ | $20.5_{\uparrow 4.3}$ | $53.6_{\uparrow 14.4}$ |

TABLE 8: Results on open-vocabulary object detection on MS COCO val set. F-ViT is the two-stage detector baseline built on the frozen original CLIP ViT, and † means the ViT backbone is initialized with the CLIP model fine-tuned (FT) with the corresponding method on the dataset in brackets (CC3M or MS COCO Karpathy trainset). Ordinary FT is equivalent to just using the global loss in the fine-tuning.

| Method | Backbone | $AP_{50}^{novel}$ | $AP_{50}^{base}$ | $AP_{50}^{all}$ |
|---|---|---|---|---|
| OV-RCNN [85] | ResNet50 | 17.5 | 41.0 | 34.9 |
| RegionCLIP [13] | ResNet50 | 26.8 | 54.8 | 47.5 |
| Detic [86] | ResNet50 | 27.8 | 51.1 | 45.0 |
| VLDet [87] | ResNet50 | 32.0 | 50.6 | 45.8 |
| F-VLM [56] | ResNet50 | 28.0 | - | 39.6 |
| CORA [57] | ResNet50 | 35.1 | 35.5 | 35.4 |
| RO-ViT [27] | ViT-B/16 | 30.2 | - | 41.5 |
| RO-ViT [27] | ViT-L/16 | 33.0 | - | 47.7 |
| F-ViT | ViT-B/16 | 19.4 | 43.3 | 37.0 |
| +CLIPSelf [28] (CC3M)† | ViT-B/16 | 13.4 | 39.3 | 32.5 |
| +Ordinary FT (CC3M)† | ViT-B/16 | 19.5 | 43.4 | 37.1 |
| +Our FT (CC3M)† | ViT-B/16 | $27.4_{\uparrow 8.0}$ | $43.8_{\uparrow 0.5}$ | $39.5_{\uparrow 2.5}$ |
| +CLIPSelf [28] (MS COCO)† | ViT-B/16 | 25.2 | 42.2 | 37.7 |
| +Ordinary FT (MS COCO)† | ViT-B/16 | 20.1 | 43.8 | 37.6 |
| +Our FT (MS COCO)† | ViT-B/16 | $26.7_{\uparrow 7.3}$ | $44.2_{\uparrow 0.9}$ | $39.6_{\uparrow 2.6}$ |
| F-ViT | ViT-L/14 | 28.3 | 52.5 | 46.2 |
| +Ordinary FT (CC3M)† | ViT-L/14 | 31.1 | 53.4 | 47.6 |
| +Our FT (CC3M)† | ViT-L/14 | $39.4_{\uparrow 11.1}$ | $53.6_{\uparrow 1.1}$ | $49.9_{\uparrow 3.7}$ |

TABLE 9: Evaluating the influence of fine-grained fine-tuning on image-level representation by a zero-shot retrieval task using Flicker30k. R@i denotes the recall accuracy with top $i$ matching. CLIP is the pre-trained CLIP model from EVA-CLIP [83], which is the base model the fine-tuning (FT) methods adopt. Ordinary FT is equivalent to just using the global loss in the fine-tuning.

| | | text-to-image | | | image-to-text | | |
| Method | Model | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 |
|---|---|---|---|---|---|---|---|
| CLIP [83] | ViT-B/16 | 73.6 | 90.9 | 94.8 | 88.6 | 97.1 | 99.1 |
| +CLIPSelf [28] | ViT-B/16 | 46.9 | 72.8 | 81.1 | 51.4 | 76.9 | 85.5 |
| +Ordinary FT | ViT-B/16 | 72.7 | 90.7 | 94.5 | 86.7 | 96.7 | 98.5 |
| +Our FT | ViT-B/16 | 72.8 | 91.0 | 94.7 | 88.3 | 97.4 | 98.9 |
| CLIP [83] | ViT-L/14 | 78.8 | 93.9 | 96.8 | 90.4 | 98.8 | 99.4 |
| +Ordinary FT | ViT-L/14 | **80.2** | 94.5 | **97.1** | 91.6 | 98.9 | 99.7 |
| +Our FT | ViT-L/14 | **80.2** | **94.6** | **97.1** | **92.1** | **99.1** | **99.8** |

the effectiveness of our FT in maintaining both the global representation and stable image-level matching.

## 7 CONCLUSION

In this paper, we propose Grad-ECLIP, a novel white-box gradient-based visual and textual explanation method for CLIP, the dual-encoder pre-trained model for image-text matching. Grad-ECLIP is applicable to both the image and text encoders, producing heat maps that indicate the importance of image regions or words for the image-text matching score. Qualitative and quantitative evaluations demonstrate the advantages of Grad-ECLIP compared with existing explanation methods designed for transformers/CLIP, and the adaptation experiments exhibit the generalizability of our method. We also adopt Grad-ECLIP to analyze the properties of the pre-trained CLIP model, where we discover its ability of concept decomposition and addibility, advantages/limitations on different attribute identification, and its bias towards learning concrete words over abstract words. By introducing these analyses as examples, we hope the proposed interpretation method can be used to help with both development and understanding of VLMs. Finally, we propose a fine-tuning scheme, which adopts the Grad-ECLIP explanation map to obtain region-text pairs for a local loss that boosts the fine-grained understanding of CLIP. In future work, we will consider how to associate individual words from the sentence to regions in the image, and vice versa. Future work can also consider how to extend Grad-ECLIP to other VLMs with modified dual-encoder architectures, e.g ALBEF with additional cross-attention layers to fuse image and text features.

the CLIP model. Finally, we conduct the experiments with the ViT-L/14 architecture and further improve the OVD performance on the novel categories by a large extent. Compared with the existing OVD methods, mostly relying on ResNet-based encoder or modified ViT encoder, and requiring pre-training from scratch on prepared large-scale data with extra region information, our method achieves superior performance with just fine-tuning the CLIP on the easily-obtained image-text pairs.

**Evaluating the influence on image-level representation.** We next explore the influence of our fine-tuning on the image-level representation ability, to see if there are negative effects to image-level representation when improving fine-grained representations. We report the recall accuracy of image-to-text and text-to-image retrieval task with the Flickr30k [33] validation set. As the results shown in Tab. 9, after fine-tuning the CLIP ViT-B/16 model on the CC3M dataset, our FT has successfully preserve the image-level retrieval performance, which is similar to the ordinary FT. In contrast, CLIPSelf has largely lost the ability of its global image-level representation. Finally, the superior results on ViT-L/14 model compared with the baseline CLIP further demonstrate

# ACKNOWLEDGEMENTS

# REFERENCES

[1] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *ICML*, 2021.

[2] S. Changpinyo, P. Sharma, N. Ding, and R. Soricut, "Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts," in *CVPR*, 2021, pp. 3558–3568.

[3] J. Cha, K. Lee, S. Park, and S. Chun, "Domain generalization by mutual-information regularization with pre-trained models," in *ECCV*, 2022.

[4] H. Luo, L. Ji, M. Zhong, Y. Chen, W. Lei, N. Duan, and T. Li, "Clip4clip: An empirical study of clip for end to end video clip retrieval and captioning," *Neurocomputing*, vol. 508, pp. 293–304, 2022.

[5] Z. Wang, Y. Lu, Q. Li, X. Tao, Y. Guo, M. Gong, and T. Liu, "Cris: Clip-driven referring image segmentation," in *CVPR*, 2022.

[6] J. Xu, S. De Mello, S. Liu, W. Byeon, T. Breuel, J. Kautz, and X. Wang, "Groupvit: Semantic segmentation emerges from text supervision," in *CVPR*, 2022.

[7] J. Yu, Z. Wang, V. Vasudevan, L. Yeung, M. Seyedhosseini, and Y. Wu, "Coca: Contrastive captioners are image-text foundation models," *arXiv:2205.01917*, 2022.

[8] J. Li, D. Li, C. Xiong, and S. Hoi, "Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation," in *ICML*, 2022, pp. 12 888–12 900.

[9] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, "Learning to prompt for vision-language models," *IJCV*, vol. 130, no. 9, pp. 2337–2348, 2022.

[10] G. Chen, W. Yao, X. Song, X. Li, Y. Rao, and K. Zhang, "Prompt learning with optimal transport for vision-language models," *arXiv:2210.01253*, 2022.

[11] X. Li, X. Yin, C. Li, P. Zhang, X. Hu, L. Zhang, L. Wang, H. Hu, L. Dong, F. Wei *et al.*, "Oscar: Object-semantics aligned pre-training for vision-language tasks," in *ECCV*, 2020.

[12] J. Wang, P. Zhou, M. Z. Shou, and S. Yan, "Position-guided text prompt for vision-language pre-training," in *CVPR*, 2023, pp. 23 242–23 251.

[13] Y. Zhong, J. Yang, P. Zhang, C. Li, N. Codella, L. H. Li, L. Zhou, X. Dai, L. Yuan, Y. Li *et al.*, "Regionclip: Region-based language-image pretraining," in *CVPR*, 2022, pp. 16 793–16 803.

[14] S. Abnar and W. Zuidema, "Quantifying attention flow in transformers," in *ACL*, 2020, pp. 4190–4197.

[15] H. Chefer, S. Gur, and L. Wolf, "Transformer interpretability beyond attention visualization," in *CVPR*, 2021.

[16] ——, "Generic attention-model explainability for interpreting bi-modal and encoder-decoder transformers," in *ICCV*, 2021, pp. 397–406.

[17] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek, "On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation," *PloS one*, vol. 10(7):e0130140, 2015.

[18] Y. Li, H. Wang, Y. Duan, and X. Li, "Clip surgery for better explainability with enhancement in open-vocabulary tasks," *arXiv:2304.05653*, 2023.

[19] C. Zhou, C. C. Loy, and B. Dai, "Extract free dense labels from clip," in *ECCV*. Springer, 2022, pp. 696–712.

[20] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *ICCV*, 2017, pp. 618–626.

[21] V. Petsiuk, A. Das, and K. Saenko, "Rise: Randomized input sampling for explanation of black-box models," *arXiv:1806.07421*, 2018.

[22] Y. Wang, T. G. Rudner, and A. G. Wilson, "Visual explanations of image-text representations via multi-modal information bottleneck attribution," *NeurIPS*, 2024.

[23] Y. Li, H. Wang, Y. Duan, H. Xu, and X. Li, "Exploring visual interpretability for contrastive language-image pre-training," *arXiv:2209.07046*, 2022.

[24] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *ECCV*, 2014, pp. 818–833.

[25] P.-T. Jiang, C.-B. Zhang, Q. Hou, M.-M. Cheng, and Y. Wei, "Layer-cam: Exploring hierarchical class activation maps for localization," *TIP*, vol. 30, pp. 5875–5888, 2021.

[26] S. Srinivas and F. Fleuret, "Full-gradient representation for neural network visualization," in *NeurIPS*, vol. 32, 2019.

[27] D. Kim, A. Angelova, and W. Kuo, "Region-aware pretraining for open-vocabulary object detection with vision transformers," in *CVPR*, 2023.

[28] S. Wu, W. Zhang, L. Xu, S. Jin, X. Li, W. Liu, and C. C. Loy, "Clipself: Vision transformer distills itself for open-vocabulary dense prediction," *arXiv preprint arXiv:2310.01403*, 2023.

[29] C. Zhao, K. Wang, X. Zeng, R. Zhao, and A. B. Chan, "Gradient-based visual explanation for transformer-based clip," in *ICML*, 2024.

[30] Z. Wang, X. Liu, H. Li, L. Sheng, J. Yan, X. Wang, and J. Shao, "Camp: Cross-modal adaptive message passing for text-image retrieval," in *ICCV*, 2019.

[31] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *ICML*, 2015, pp. 2048–2057.

[32] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, "Vqa: Visual question answering," in *ICCV*, 2015.

[33] B. A. Plummer, L. Wang, C. M. Cervantes, J. C. Caicedo, J. Hockenmaier, and S. Lazebnik, "Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models," in *ICCV*, 2015.

[34] A. Chattopadhay, A. Sarkar, P. Howlader, and V. N. Balasubramanian, "Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks," in *WACV*, 2018, pp. 839–847.

[35] H. G. Ramaswamy *et al.*, "Ablation-cam: Visual explanations for deep convolutional network via gradient-free localization," in *WACV*, 2020.

[36] H. Wang, Z. Wang, M. Du, F. Yang, Z. Zhang, S. Ding, P. Mardziel, and X. Hu, "Score-cam: Score-weighted visual explanations for convolutional neural networks," in *CVPR Workshops*, 2020, pp. 24–25.

[37] H. Wang, R. Naidu, J. Michael, and S. S. Kundu, "Ss-cam: Smoothed score-cam for sharper visual feature localization," *arXiv:2006.14255*, 2020.

[38] M. T. Ribeiro, S. Singh, and C. Guestrin, "" why should i trust you?" explaining the predictions of any classifier," in *ACM SIGKDD*, 2016.

[39] R. C. Fong and A. Vedaldi, "Interpretable explanations of black boxes by meaningful perturbation," in *ICCV*, 2017, pp. 3429–3437.

[40] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," *NeurIPS*, vol. 30, 2017.

[41] J. Wagner, J. M. Kohler, T. Gindele, L. Hetzel, J. T. Wiedemer, and S. Behnke, "Interpretable and fine-grained visual explanations for convolutional neural networks," in *CVPR*, 2019, pp. 9097–9107.

[42] J. Lee, J. Yi, C. Shin, and S. Yoon, "Bbam: Bounding box attribution map for weakly supervised semantic and instance segmentation," in *CVPR*, 2021.

[43] V. Petsiuk, R. Jain, V. Manjunatha, V. I. Morariu, A. Mehra, V. Ordonez, and K. Saenko, "Black-box explanation of object detectors via saliency maps," in *CVPR*, 2021, pp. 11 443–11 452.

[44] G. Montavon, S. Lapuschkin, A. Binder, W. Samek, and K.-R. Müller, "Explaining nonlinear classification decisions with deep taylor decomposition," *Pattern recognition*, vol. 65, pp. 211–222, 2017.

[45] W.-J. Nam, S. Gur, J. Choi, L. Wolf, and S.-W. Lee, "Relative attributing propagation: Interpreting the comparative contributions of individual units in deep neural networks," in *AAAI*, 2020.

[46] A. Shrikumar, P. Greenside, and A. Kundaje, "Learning important features through propagating activation differences," in *ICML*, 2017.

[47] J. Gu, Y. Yang, and V. Tresp, "Understanding individual decisions of cnns via contrastive backpropagation," in *ACCV*, 2019, pp. 119–134.

[48] B. K. Iwana, R. Kuroki, and S. Uchida, "Explaining convolutional neural networks using softmax gradient layer-wise relevance propagation," in *ICCVW*, 2019.

[49] Y. Qiang, D. Pan, C. Li, X. Li, R. Jang, and D. Zhu, "Attcat: Explaining transformers via attentive class activation tokens," *NeurIPS*, 2022.

[50] W. Xie, X.-H. Li, C. C. Cao, and N. L. Zhang, "Vit-cx: Causal explanation of vision transformers," pp. 1569–1577, 2023.

[51] L. Yu and W. Xiang, "X-pruner: explainable pruning for vision transformers," in *CVPR*, 2023, pp. 24 355–24 363.

[52] M. Xu, Z. Zhang, F. Wei, Y. Lin, Y. Cao, H. Hu, and X. Bai, "A simple baseline for zeroshot semantic segmentation with pre-trained vision-language model," *ECCV*, 2022.

[53] F. Liang, B. Wu, X. Dai, K. Li, Y. Zhao, H. Zhang, P. Zhang, P. Vajda, and D. Marculescu, "Open-vocabulary semantic segmentation with mask-adapted clip," in *CVPR*, 2023, pp. 7061–7070.

[54] X. Gu, T.-Y. Lin, W. Kuo, and Y. Cui, "Open-vocabulary object detection via vision and language knowledge distillation," *ICLR*, 2022.

[55] Y. Du, F. Wei, Z. Zhang, M. Shi, Y. Gao, and G. Li, "Learning to prompt for open-vocabulary object detection with vision-language model," in *CVPR*, 2022, pp. 14 084–14 093.

[56] W. Kuo, Y. Cui, X. Gu, A. Piergiovanni, and A. Angelova, "F-vlm: Open-vocabulary object detection upon frozen vision and language models," in *ICLR*, 2023.

[57] X. Wu, F. Zhu, R. Zhao, and H. Li, "Cora: Adapting clip for open-vocabulary detection with region prompting and anchor pre-matching,"

in *CVPR*, 2023, pp. 7031–7040.

[58] Q. Yu, J. He, X. Deng, X. Shen, and L.-C. Chen, "Convolutions die hard: Open-vocabulary segmentation with single frozen convolutional clip," *NeurIPS*, vol. 36, 2024.

[59] L. H. Li, P. Zhang, H. Zhang, J. Yang, C. Li, Y. Zhong, L. Wang, L. Yuan, L. Zhang, J.-N. Hwang *et al.*, "Grounded language-image pre-training," in *CVPR*, 2022, pp. 10 965–10 975.

[60] Y.-C. Chen, L. Li, L. Yu, A. El Kholy, F. Ahmed, Z. Gan, Y. Cheng, and J. Liu, "Uniter: Universal image-text representation learning," in *ECCV*, 2020.

[61] P. Zhang, X. Li, X. Hu, J. Yang, L. Zhang, L. Wang, Y. Choi, and J. Gao, "Vinvl: Revisiting visual representations in vision-language models," in *CVPR*, 2021, pp. 5579–5588.

[62] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu *et al.*, "Grounding dino: Marrying dino with grounded pre-training for open-set object detection," *arXiv preprint arXiv:2303.05499*, 2023.

[63] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma *et al.*, "Visual genome: Connecting language and vision using crowdsourced dense image annotations," *IJCV*, vol. 123, pp. 32–73, 2017.

[64] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *NeurIPS*, vol. 28, 2015.

[65] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.

[66] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *IJCV*, vol. 115, pp. 211–252, 2015.

[67] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *ECCV*, 2014.

[68] D. Hendrycks, S. Basart, N. Mu, S. Kadavath, F. Wang, E. Dorundo, R. Desai, T. Zhu, S. Parajuli, M. Guo *et al.*, "The many faces of robustness: A critical analysis of out-of-distribution generalization," in *ICCV*, 2021, pp. 8340–8349.

[69] H. Wang, S. Ge, Z. Lipton, and E. P. Xing, "Learning robust global representations by penalizing local predictive power," *NeurIPS*, 2019.

[70] D. Hendrycks, K. Zhao, S. Basart, J. Steinhardt, and D. Song, "Natural adversarial examples," in *CVPR*, 2021, pp. 15 262–15 271.

[71] P. Sharma, N. Ding, S. Goodman, and R. Soricut, "Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning," in *ACL*, 2018, pp. 2556–2565.

[72] B. Boecking, N. Usuyama, S. Bannur, D. C. Castro, A. Schwaighofer, S. Hyland, M. Wetscherek, T. Naumann, A. Nori, J. Alvarez-Valle *et al.*, "Making the most of text semantics to improve biomedical vision–language processing," in *ECCV*, 2022.

[73] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[74] J. Li, R. Selvaraju, A. Gotmare, S. Joty, C. Xiong, and S. C. H. Hoi, "Align before fuse: Vision and language representation learning with momentum distillation," *NeurIPS*, vol. 34, pp. 9694–9705, 2021.

[75] W. Samek, A. Binder, G. Montavon, S. Lapuschkin, and K.-R. Müller, "Evaluating the visualization of what a deep neural network has learned," *IEEE Trans. Neural Netw Learn Syst*, vol. 28(11), pp. 2660–73, 2016.

[76] J. Zhang, S. A. Bargal, Z. Lin, J. Brandt, X. Shen, and S. Sclaroff, "Top-down neural attention by excitation backprop," *IJCV*, vol. 126(10), pp. 1084–102, 2018.

[77] C. Zhao and A. B. Chan, "Odam: Gradient-based instance-specific visual explanation for object detection," in *ICLR*, 2022.

[78] S. Gao, Z.-Y. Li, M.-H. Yang, M.-M. Cheng, J. Han, and P. Torr, "Large-scale unsupervised semantic segmentation," *TPAMI*, 2022.

[79] J. Johnson, B. Hariharan, L. Van Der Maaten, L. Fei-Fei, C. Lawrence Zitnick, and R. Girshick, "Clevr: A diagnostic dataset for compositional language and elementary visual reasoning," in *CVPR*, 2017.

[80] M. Brysbaert, A. B. Warriner, and V. Kuperman, "Concreteness ratings for 40 thousand generally known english word lemmas," *Behavior research methods*, vol. 46, pp. 904–911, 2014.

[81] A. Gokaslan and V. Cohen, "Openwebtext corpus," http://Skylion007.github.io/OpenWebTextCorpus, 2019.

[82] S. Bird, E. Klein, and E. Loper, *Natural language processing with Python: analyzing text with the natural language toolkit.* O'Reilly Media, Inc., 2009.

[83] Q. Sun, Y. Fang, L. Wu, X. Wang, and Y. Cao, "Eva-clip: Improved training techniques for clip at scale," *preprint arXiv:2303.15389*, 2023.

[84] A. Kirillov, K. He, R. Girshick, C. Rother, and P. Dollár, "Panoptic segmentation," in *CVPR*, 2019, pp. 9404–9413.

[85] A. Zareian, K. D. Rosa, D. H. Hu, and S.-F. Chang, "Open-vocabulary object detection using captions," in *CVPR*, 2021, pp. 14 393–14 402.

[86] X. Zhou, R. Girdhar, A. Joulin, P. Krähenbühl, and I. Misra, "Detecting twenty-thousand classes using image-level supervision," in *ECCV*, 2022.

[87] C. Lin, P. Sun, Y. Jiang, P. Luo, L. Qu, G. Haffari, Z. Yuan, and J. Cai, "Learning object-language alignments for open-vocabulary object detection," *arXiv preprint arXiv:2211.14843*, 2022.

[88] X. Chen, H. Fang, T.-Y. Lin, R. Vedantam, S. Gupta, P. Dollár, and C. L. Zitnick, "Microsoft coco captions: Data collection and evaluation server," *arXiv preprint arXiv:1504.00325*, 2015.

**Chenyang Zhao** received the B.Eng. degree in Electrical Engineering from Xiamen University, Xiamen, China, and M.S. degree in Computer Science from School of Electronic and Computer Engineering, Peking University, Shenzhen, China, in 2016 and 2019, respectively. She is currently working towards the Ph.D. degree in Computer Science at the City University of Hong Kong. Her research interests include explainable AI and object detection.

**Kun Wang** is a senior researcher at SenseTime Group Limited. He holds an MPhil degree from the Department of Electronic Engineering at the Chinese University of Hong Kong. His research interests focus on computer vision and representation learning. Presently, he is engaged in developing applications utilizing large language models and large multi-modal models within the industry.

**Janet H. Hsiao** received the B.S. degree in Computer Science & Information Engineering from National Taiwan University, the M.S. degree in Computing Science from Simon Fraser University, and the Ph.D. degree in Informatics from University of Edinburgh. She is currently a Professor in the Division of Social Science and Department of Computer Science & Engineering at Hong Kong University of Science and Technology. She is also a Fellow of the Cognitive Science Society and serves on the Governing Board. Her research interests include cognitive science, computational modelling, learning and visual cognition, and explainable AI.

**Antoni B. Chan** received the B.S. and M.Eng. degrees in electrical engineering from Cornell University, Ithaca, NY, in 2000 and 2001, and the Ph.D. degree in electrical and computer engineering from the University of California, San Diego (UCSD), San Diego, in 2008. He is currently a Professor in the Department of Computer Science, City University of Hong Kong. His research interests include computer vision, machine learning, pattern recognition, and music analysis.