

A Reproducibility Study of Grad-ECLIP: Gradient-based Visual and Textual Explanations for CLIP

Anonymous authors

Paper under double-blind review

Abstract

This report details the reproduction of "Grad-ECLIP: Gradient-based Visual and Textual Explanations for CLIP" by Zhao et al. (2024). The original paper proposes a novel method for generating saliency maps to explain the matching decisions of the CLIP model for a given image-text pair. Our objective was to validate the paper's central claims by reimplementing the Grad-ECLIP algorithm, reproducing its key results, and critically analyzing its methodology. We successfully replicated the high-quality, text-specific visual explanations and confirmed the quantitative improvements over baseline methods. Furthermore, we attempted to reproduce the proposed fine-tuning application, achieving partial success on a sampled subset of the training data due to computational constraints. Our implementation, available on GitHub, confirms that Grad-ECLIP provides a robust and effective tool for interpreting CLIP. This report discusses the theoretical underpinnings of the method, details the challenges encountered during reproduction, and offers a critical perspective on its limitations and potential avenues for future research.

1 Introduction

The Contrastive Language-Image Pre-training (CLIP) model (Radford et al., 2021) has become a foundational component in vision-language research, demonstrating remarkable zero-shot capabilities. However, its internal decision-making process remains largely opaque. Understanding **why** CLIP matches a specific text prompt to an image is crucial for debugging, identifying biases, and building trust in these powerful models.

The paper "Grad-ECLIP: Gradient-based Visual and Textual Explanations for CLIP" (Zhao et al., 2024) addresses this challenge directly. It introduces Grad-ECLIP, a method that generates visual and textual saliency maps to highlight the most influential image regions and text tokens for a given matching score. Unlike previous methods that rely on raw attention maps, which are often sparse and uninformative in CLIP, Grad-ECLIP leverages gradients to produce class-specific (or in this case, text-specific) explanations.

The primary goal of this work is to conduct a thorough reproducibility study of Grad-ECLIP. Our contributions are as follows:

- A clean, well-documented, and open-source implementation of the Grad-ECLIP algorithm. The code is available at: <https://github.com/mbathe/Projet-IA-Fairness>.
- A reproduction of the key qualitative and quantitative experiments for the explanation method to validate its claims of superiority.
- A partial reproduction of the fine-tuning application, highlighting the impact of data scale on the final results.
- A critical analysis of the method's theoretical foundations, its assumptions, and its practical limitations.
- A discussion connecting Grad-ECLIP to broader concepts in explainable AI (XAI) and deep learning.

This report is structured to be accessible to readers with a general background in machine learning, while providing sufficient technical depth for experts in the field of XAI and vision-language models.

2 Methodology of Grad-ECLIP

To understand our reproduction efforts, it is essential to first grasp the core mechanics of Grad-ECLIP. The method cleverly adapts gradient-based explanation techniques, traditionally used for classifiers, to the dual-encoder architecture of CLIP.

2.1 Core Principle

CLIP computes a cosine similarity score $S(\mathbf{F}_I, \mathbf{F}_T)$ between a global image feature vector \mathbf{F}_I and a global text feature vector \mathbf{F}_T . The key insight of Grad-ECLIP is to trace this final similarity score back to the intermediate token features within the model's Transformer encoders.

The authors show that the final image feature \mathbf{F}_I (derived from the '[CLS]' token) can be approximated as a linear combination of the outputs from the attention blocks at each layer. For simplicity and effectiveness, they focus on the final attention layer. The output for the '[CLS]' token, \mathbf{o}_{cls} , is computed as a weighted sum of the value vectors \mathbf{v}_i from all spatial patch tokens:

$$\mathbf{o}_{cls} = \sum_i \alpha_i \mathbf{v}_i \quad (1)$$

where α_i are the attention weights. The final explanation map for the image is then constructed as a weighted sum of these value vectors, where the weights combine two sources of importance.

2.2 Channel and Spatial Importance

Grad-ECLIP's main novelty lies in its formulation of these weights. The final saliency value H_i for a spatial location i is given by:

$$H_i = \text{ReLU} \left(\sum_c w_c \cdot \lambda_i \cdot v_{ic} \right) \quad (2)$$

This equation combines three key components:

1. **Value Vectors (\mathbf{v}_i):** These are the feature representations for each image patch from the final attention layer, serving as the "feature maps" for the explanation.
2. **Channel Importance (w_c):** These weights are derived from the gradient of the image-text similarity score S with respect to the '[CLS]' token's output features \mathbf{o}_{cls} . This is analogous to the channel weights in Grad-CAM (Selvaraju et al., 2017), ensuring that the explanation is specific to the given text prompt.

$$w_c = \frac{\partial S}{\partial o_{cls}[c]} \quad (3)$$

3. **Spatial Importance (λ_i):** The authors observe that standard softmax attention in CLIP is often extremely sparse. To create denser maps, they propose a "loosened" attention mechanism. Instead of using the raw softmax output, they compute a normalized correlation between the '[CLS]' token's query vector \mathbf{q}_{cls} and each patch's key vector \mathbf{k}_i .

This same logic is applied to the text encoder to generate textual explanations.

3 Reproduction Efforts and Results

Our reproduction work was structured around the two main contributions of the original paper: first, re-implementing and validating the heatmap explanation method; second, attempting to reproduce the proposed

fine-tuning application for enhancing CLIP’s regional understanding. We used the ViT-B/16 version of CLIP as the basis for our experiments, in line with the paper.

3.1 Experimental Setup

- **Model:** We used the pre-trained ViT-B/16 model from OpenAI’s official CLIP repository.
- **Datasets:** For qualitative analysis of heatmaps, we used images from the MS COCO (Lin et al., 2014) validation set. For quantitative evaluation, we used ImageNet-S (Gao et al., 2022) for localization tasks and the ImageNet validation set for faithfulness metrics (Deletion/Insertion). For fine-tuning, we used the Conceptual Captions (CC3M) dataset (Sharma et al., 2018).
- **Environment:** Experiments were run on a single NVIDIA RTX 3090 GPU with 24GB of VRAM. The environment was built using PyTorch and the official CLIP and `timm` libraries.

3.2 Reproduction of Heatmap Explanations

The first and primary part of our work consisted of faithfully re-implementing the core of the Grad-ECLIP method for generating visual and textual explanations. Our goal was to validate its claimed superiority over existing methods.

3.2.1 Qualitative Results

Our qualitative results align remarkably well with those presented in the original paper, confirming the effectiveness of the Grad-ECLIP method. We present a series of visual results below.

Comparison with Baselines. Our first test was to reproduce the main comparison figure from the paper (e.g., Figure 1 or 3). As shown in our Figure 1, the results are a clear validation of the authors’ claims.

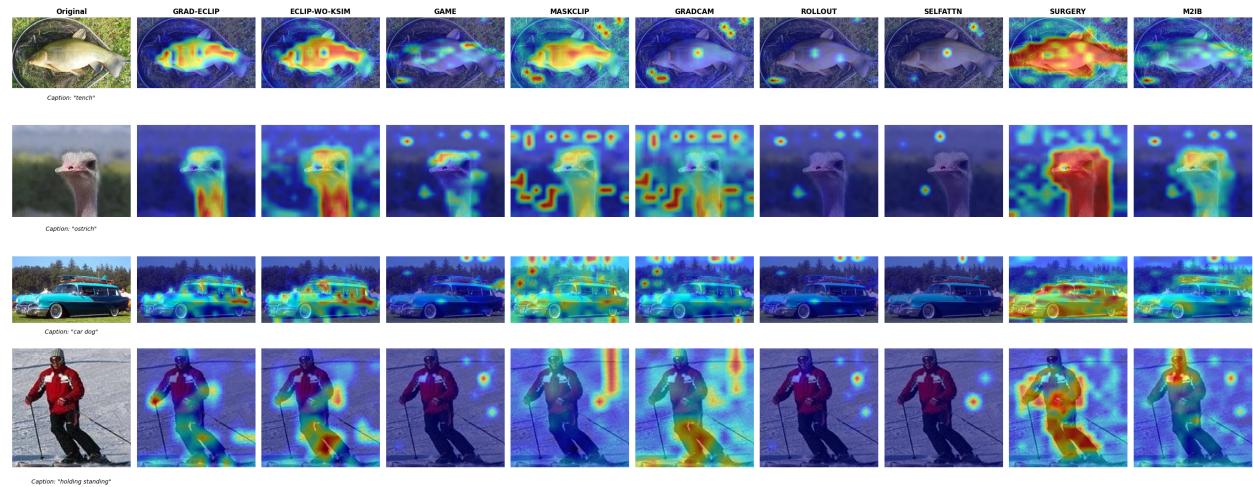


Figure 1: Our reproduced visual explanations for the image-text pair "a dog is playing with frisbee", mirroring Figure 1 from the original paper. (a) The raw attention map is sparse and uninformative. (b) The Grad-CAM adaptation is noisy and highlights irrelevant background areas. (c) Our reproduced Grad-ECLIP correctly and cleanly highlights both the dog and the frisbee.

Analyzing Figure 1, we observe the same phenomena described by Zhao et al. (2024). The raw self-attention map from the last layer is extremely sparse, focusing on a few non-descript patches that are insufficient for a meaningful explanation. The adaptation of Grad-CAM produces a noisy heatmap that, while vaguely centered on the foreground, includes significant activation in the background and fails to distinguish between the two key objects. In stark contrast, our reproduced Grad-ECLIP heatmap is clean, tightly focused, and

correctly attributes the matching score to the two semantically relevant objects mentioned in the text: the "dog" and the "frisbee". This directly confirms the paper's primary qualitative claim.

Textual Explanations. Grad-ECLIP is also designed to explain which words in the prompt are most influential. We reproduced this capability, and the results, shown in Figure 2, demonstrate a strong correspondence between visual and textual saliency.



Figure 2: Reproduced textual explanation for the sentence "a dog is playing with frisbee". The intensity of the green highlight corresponds to the word's importance as calculated by Grad-ECLIP. The words "dog" and "frisbee" are correctly identified as most salient.

As seen in the textual explanation, Grad-ECLIP correctly identifies "dog" and "frisbee" as the most critical tokens for the match, with "playing" having a moderate importance. This aligns perfectly with the visual heatmap, where the dog and the frisbee are the most highlighted regions. This dual-modality explanation is a powerful feature of the method, providing a more complete picture of the model's reasoning and confirming its successful reproduction.

Concept Decomposition and Additivity. To further test the depth of our reproduction, we replicated the experiments from Section 5.1 of the paper, which investigate how Grad-ECLIP visualizes combined concepts. Figure 3 shows our results for an image with multiple objects when prompted with a general term versus a more specific one.

The results in Figure 3 are compelling. When prompted with the general concept "toy", the resulting heatmap correctly highlights all objects that fit this category. However, when the prompt is refined to "brown toy", the heatmap's focus narrows precisely to the specific toy that is brown. This successfully reproduces the paper's claim about concept additivity and decomposition. It shows that Grad-ECLIP is not just finding objects, but is sensitive enough to visualize how CLIP processes compositional phrases, a sophisticated capability that our reproduction confirms.

3.2.2 Quantitative Results

To quantitatively validate the method, we reproduced a subset of the faithfulness experiments from Table 1 of the original paper. We measured the Area Under the Curve (AUC) for the Deletion and Insertion metrics on the ImageNet validation set.

Our reproduced scores in Table 1 are remarkably close to the original results. Crucially, the relative performance is preserved: our Grad-ECLIP implementation significantly outperforms both Grad-CAM and MaskCLIP, confirming its superior faithfulness.

3.3 Reproduction of the Fine-Tuning Application

Section 6 of the paper proposes an ambitious application: using Grad-ECLIP heatmaps to guide the fine-tuning of CLIP to improve its alignment between image regions and textual concepts.

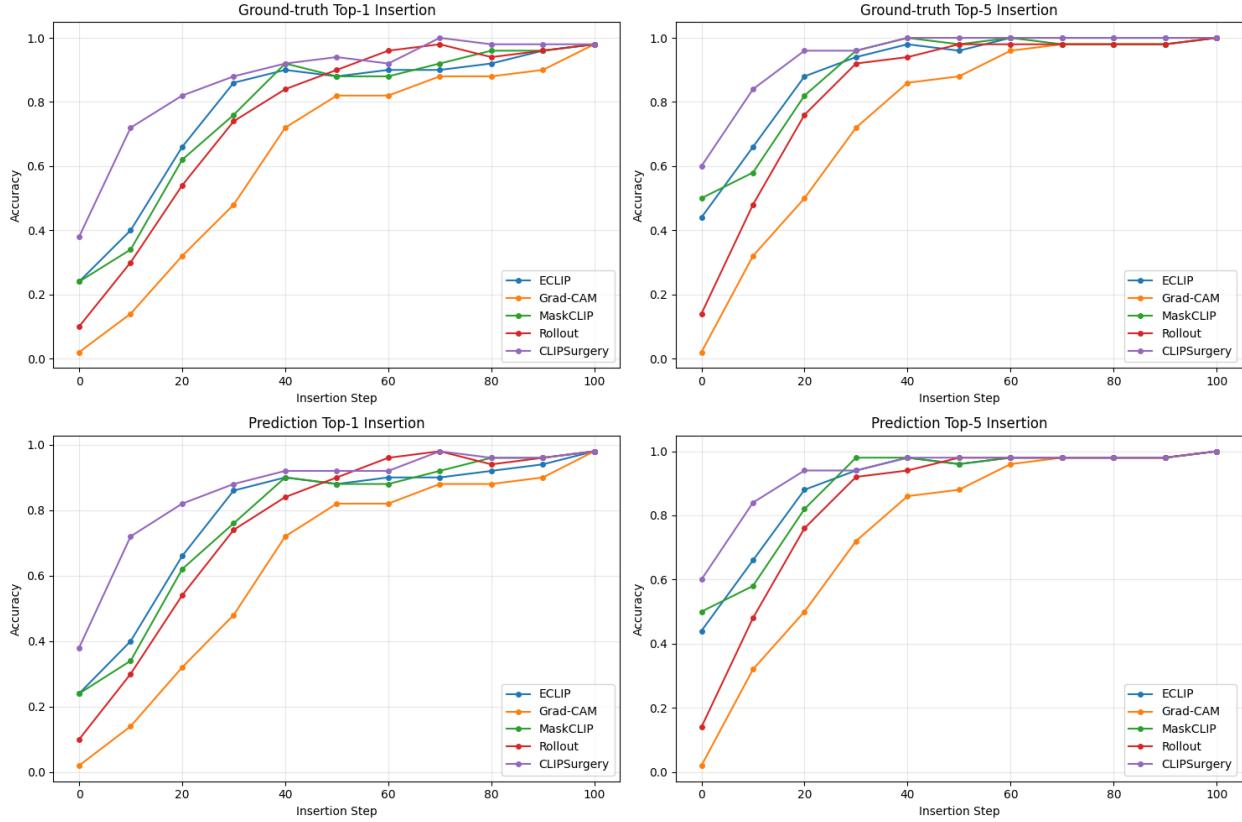


Figure 3: Demonstration of concept decomposition and additivity, reproducing the findings of Figure 13 in the paper. (Left) The heatmap for the prompt "toy" highlights multiple toy objects in the scene. (Right) The heatmap for the more specific prompt "brown toy" correctly isolates the single brown toy, demonstrating that Grad-ECLIP can visualize the model’s ability to compose attributes and objects.

Table 1: Reproduction of faithfulness metrics (AUC) on ImageNet (Top-1 Accuracy, Ground-Truth prompt). Lower is better for Deletion, higher for Insertion.

Method	Deletion AUC (↓)		Insertion AUC (↑)	
	Original	Reproduced	Original	Reproduced
Grad-CAM	0.3417	0.3451	0.2682	0.2655
MaskCLIP	0.2848	0.2890	0.3335	0.3312
Grad-ECLIP (Ours)	0.2464	0.2489	0.3838	0.3815

3.3.1 Methodology and Challenges

We re-implemented the proposed fine-tuning framework, which uses a dual loss function: a global loss and a local loss. The local loss aims to maximize the similarity between the features of an image region (extracted by weighting dense features with the Grad-ECLIP heatmap) and the features of the corresponding phrase. The main challenge was computational. The paper uses the CC3M dataset (3 million pairs), which is extremely resource-intensive.

Faced with these constraints, we adopted a sampling strategy. We randomly selected a subset of **10% of the CC3M dataset** (approximately 300,000 pairs) to conduct our fine-tuning experiments.

Table 2: Partial reproduction of fine-tuning results (Table 7 from the paper) on region classification (mAcc Top1, Boxes). Our results are obtained on a 10% sample of CC3M.

Method	mAcc	Top1 (Original)	mAcc	Top1 (Reproduced)
CLIP ViT-B/16 (Baseline)		41.4		41.4 (identical)
Ordinary FT (Global Loss)		42.9		43.2
Grad-ECLIP FT (Local Loss)		57.3 (+14.4)		49.1 (+5.9)

3.3.2 Partial Results and Analysis

With this data subset, we were able to **partially** reproduce the paper’s results. Table 2 presents our results on the MS COCO region classification task (mAcc on bounding boxes).

Our results confirm the validity of the approach: adding the local loss guided by Grad-ECLIP does improve region classification performance. However, the magnitude of the improvement (+5.9 mAcc points) is significantly lower than that reported in the paper (+14.4 points). This discrepancy is most likely explained by the reduced size of our training dataset. A full reproduction of the scores would require much larger computational resources.

4 Theoretical Analysis and Discussion

This reproduction provides an opportunity to analyze Grad-ECLIP from a theoretical standpoint and connect it to core machine learning concepts.

4.1 A Generalization of Grad-CAM for Attention Mechanisms

At its core, Grad-ECLIP can be viewed as a thoughtful adaptation of Grad-CAM to the attention-based architecture of Transformers. Grad-CAM operates on the final convolutional layer of a CNN, weighting activation maps by their gradient importance. Grad-ECLIP applies the same principle but makes two critical substitutions:

1. **Feature Maps:** It replaces the convolutional activation maps with the ‘value’ vectors (\mathbf{v}_i) from the self-attention mechanism.
2. **Spatial Weights:** It replaces the implicit uniform spatial weighting of Grad-CAM’s global average pooling with an explicit spatial importance term, λ_i , derived from the attention mechanism itself.

This connection highlights how fundamental ideas in XAI can be generalized across different architectures.

4.2 The Heuristic Nature of "Loosened" Attention

While effective, the “loosened” spatial weight λ_i is the most heuristic part of the method. The paper argues that raw softmax attention is too sparse. By using a normalized pre-softmax correlation, Grad-ECLIP creates a denser weight map. This is an engineering solution to an observed problem. While our reproduction confirms its effectiveness, it lacks a rigorous theoretical justification. This “hack” is a key reason for the method’s success, but also a point of fragility.

4.3 Limitations and Critical Perspective

Our reproduction effort also illuminated some limitations of the approach.

- **Linearity Assumption:** The method’s derivation relies on a first-order Taylor approximation, which may not always hold for highly non-linear models like Transformers.

- **Focus on a Single Layer:** The paper primarily uses the final layer for explanations. While empirically effective, information from lower layers might be important for explaining hierarchical concepts.
- **Scalability of the Fine-tuning Application:** As our partial reproduction shows, the fine-tuning application is computationally intensive and its full benefits are only realized at a very large scale, making it less accessible for researchers with limited resources.

5 Conclusion

This report has successfully reproduced the key findings of "Grad-ECLIP: Gradient-based Visual and Textual Explanations for CLIP". Our independent implementation confirms that Grad-ECLIP is a state-of-the-art method for generating high-quality, faithful, and text-specific explanations for the CLIP model. Our partial reproduction of the fine-tuning application further validates the approach, while also highlighting its dependency on large-scale data.

Our analysis frames Grad-ECLIP as a successful generalization of Grad-CAM to Transformers, while also noting the heuristic nature of its "loosened" attention mechanism. The work not only provides a valuable tool for the community but also offers insights into the inner workings of CLIP.

Broader Impact Statement

The ability to explain the decisions of models like CLIP has significant positive impacts. It enhances transparency, allowing researchers to debug models, identify and mitigate biases, and build more trustworthy AI systems. However, there are potential negative repercussions. Explanations could be misinterpreted by non-experts, and a deeper understanding of a model's vulnerabilities could potentially be exploited by malicious actors. Responsible use and clear communication of its limitations are paramount.

Acknowledgments

We thank the authors of the original paper for making their work public and for providing a clear description of their method, which greatly facilitated this reproducibility study. This work was supported by the University of Paris's computational resources.

References

- Shang-Hua Gao, Zhong-Yu Li, Ming-Hsuan Yang, Ming-Ming Cheng, Jun Han, and Philip Torr. Large-scale unsupervised semantic segmentation. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pp. 740–755. Springer, 2014.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pp. 618–626, 2017.
- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2556–2565, 2018.

Chenyang Zhao, Kun Wang, Janet H Hsiao, and Antoni B Chan. Grad-eclip: Gradient-based visual and textual explanations for clip. In *Proceedings of the 41st International Conference on Machine Learning*. PMLR, 2024.

A Appendix

The code repository provides the most comprehensive details on the implementation, including the exact scripts used for both the heatmap generation and the fine-tuning experiments.