
LRP-CLIP: A Zero Shot Approach for the Explanation of the Cognitive Functions of Vision Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Cognitive explanation of deep neural networks is enabled by the use of Layer-wise
2 Relevance Propagation (LRP) heatmaps. Heatmaps indicate where a region of
3 high significance is, but current approaches rely on further computational training
4 overhead to assign semantic meaning. We propose LRP-CLIP, a zero-shot
5 framework that grounds these heatmaps in natural language using pretrained CLIP
6 (Contrastive Language-Image Pre-Training) and thereby allowing cognitive ex-
7 planations of deep learning models without further training. Our method extracts
8 relevance-based image crops and matches them against domain-specific textual
9 attributes, producing human-readable explanations without any additional training
10 or supervision. Applied to bird classification with a Vision Transformer (ViT), we
11 find that the model relies heavily on head features, while occasionally exploiting
12 spurious background cues. This dual pattern reflects both systematic heuristics
13 and potential biases in the model’s cognitive strategy. Quantitative evaluation
14 shows that our pipeline reliably localizes annotated parts and assigns semantically
15 valid labels, clearly outperforming random baselines. These results demonstrate an
16 extensible zero-shot approach to model cognition.

17 1 Motivation

18 Deep learning models achieve impressive results, but their decision-making remains opaque. These
19 black-box systems can make accurate predictions, yet offer little insight into why they arrive at a
20 certain conclusion. Layer-wise Relevance Propagation (LRP) highlights which parts of the input
21 contribute most to the output, but the resulting maps are often abstract: they show where the model
22 focuses, not what it sees or why those regions matter. To make these explanations more meaningful,
23 we propose using external knowledge to give structure to the spatial relevance information. Our
24 approach consists two steps: (1) decompose relevance hotspots into distinct sub-regions, and (2)
25 reclassify these components and translate them into natural language, so a reader understands both
26 what the model is looking at and its semantic meaning. Building on this idea, we propose LRP-CLIP,
27 a novel framework that extends LRP-based explanations by semantically interpreting the highlighted
28 sub-regions using CLIP (Contrastive Language-Image Pre-Training). CLIP aligns images and text in
29 a shared embedding space for zero-shot classification. We extract relevant input regions using LRP
30 and then classify these regions with CLIP, using domain-specific textual labels and attributes. This
31 approach produces human-readable explanations without requiring specialized expert knowledge to
32 interpret LRP maps or additional training, since it leverages CLIP directly. This approach bridges
33 low-level model reasoning with natural language, offering a modular pipeline for explainability. We
34 demonstrate the method’s feasibility on the CUB-200 bird classification dataset Wah et al. [2011]
35 using a Vision Transformer (ViT) and discuss its strengths, limitations, and future applicability to
36 other domains, including time-series data.

37 **2 Related Work**

38 The complexity of modern models, particularly Transformers, has amplified interest in explainable
 39 AI (XAI), which seeks to make decisions more transparent and trustworthy, especially in high-stakes
 40 domains.
 41 Explainability approaches fall into three categories: data-level, model-level, and post-hoc Sun et al.
 42 [2024]. Data-level methods jointly model multiple modalities to produce both predictions and ex-
 43 planations. Model-level methods use inherently interpretable architectures. Post-hoc methods, in
 44 contrast, analyze a model’s outputs in relation to its inputs, independent of training.
 45 Post-hoc methods include signal-based techniques such as LRP Ahtibat et al. [2024] and Grad-
 46 CAM Selvaraju et al. [2019], surrogate models such as LIME Ribeiro et al. [2016], and hybrids that
 47 combine them Dhore et al. [2024].
 48 Yet even high-quality relevance maps often remain inaccessible to non-experts Bove et al. [2022],
 49 highlighting the semantic gap between what the model and what humans interpret.
 50 Recent work has turned to zero-shot models such as CLIP Radford et al. [2021], which align vision
 51 and language for label-free classification using natural language prompts.
 52 Some studies use CLIP directly for explanations, generating zero-shot semantic descriptions Sammani
 53 and Deligiannis [2024].
 54 We propose a pipeline that integrates LRP for visual relevance extraction with CLIP for semantic
 55 labeling of regions, bridging the gap between faithful explanations and human understanding without
 56 retraining and limited architecture complexity.

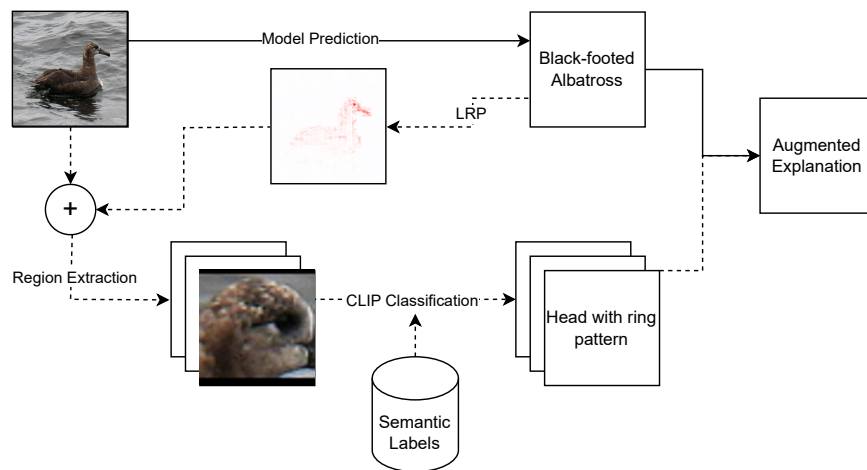


Figure 1: The general pipeline flow design of LRP-CLIP.

57 **3 Architecture of the Pipeline**

58 We propose a model agnostic modular architecture for generating post-hoc, human-interpretable
 59 explanations of classification decisions on image data (see Figure 1). The system consists of two main
 60 components: the Model Prediction Module and the Explanation Module. The Prediction Module
 61 carries out the core classification task. It receives an image as input and outputs a predicted class
 62 label. In Figure 1, the Prediction Module is depicted in the top part, where the image-based input is
 63 fed into a model that performs a prediction. In this example, the output is the class label *Black-footed*
 64 *Albatross*. The Explanation Module operates after a prediction has been made for a given input. In
 65 our approach, we use LRP to generate relevance heatmaps. These heatmaps, in combination with
 66 the original input image, are used to extract the most relevant regions. As a result, we obtain one or
 67 more image regions (from 1 to N), each of which is passed to CLIP, a vision-language model that has
 68 been primed with external knowledge (i.e. Semantic Labels). In the example shown in Figure 1, the
 69 external knowledge consists of general information about bird species.

70 Implementation

71 We fine-tuned a ViT-L (ImageNet-1000 pretrained) on the CUB-200 dataset, reaching 87.2% validation accuracy. The Explanation Module is implemented using the following methods:

73 The projection from raw input to relevant information is performed using LRP. We used LRP with ViT-specific settings Achtibat et al. [2024], enforcing positive relevance for more focused heatmaps Anders et al. [2021]. Regions of highest relevance were extracted, merged if adjacent, or split if complex (e.g., head vs. eye) based on high relevance regions, producing the largest meaningful crops. These were ordered by relevance and passed to CLIP. We used CLIP ViT-H/14 (LAION-2B pretrained) for zero-shot classification, comparing crops against CUB-200 attribute labels in text space Schuhmann et al. [2022], Radford et al. [2021], Ilharco et al. [2021].

80 4 Evaluation

81 We evaluate on the CUB-200 dataset Wah et al. [2011], which includes 12k images across 200 bird species with fine-grained annotations (labels, parts, and attributes). These detailed annotations make it well-suited for assessing both prediction accuracy and the semantic quality of explanations.

84 Findings

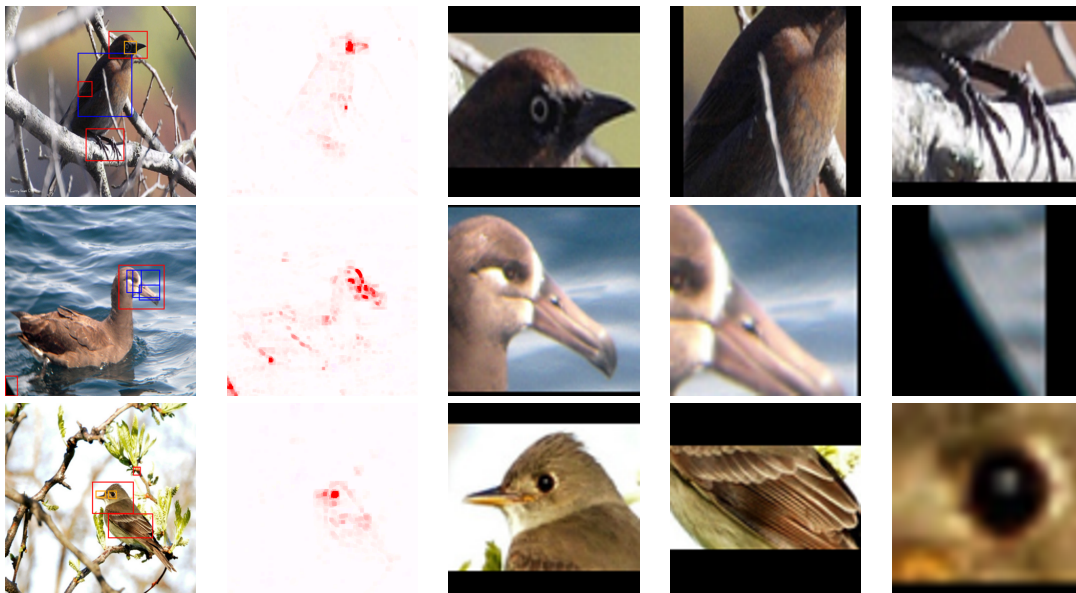


Figure 2: The three largest relevance-based crops are shown. Classifications include a rusty blackbird, a black-footed albatross, and a misclassification of an Acadian Flycatcher (predicted as Yeast Flycatcher, with Acadian Flycatcher second). The top CLIP labels of the nine crops are: (1) black forehead, (2) black back, (3) gray legs, (4) white forehead, (5) hooked bill, buff-colored, similar in size to head, (6) bird of medium (9–16 in) size, (7) rufous forehead, (8) broad buff-colored wings, (9) buff-colored eye.

85 Exemplary outputs of the pipeline (Fig. 2) show that the ViT consistently identifies the head as the most relevant region, while other highlighted areas vary by image (e.g., body, wings, or legs). CLIP generally assigns meaningful labels to these regions, though it occasionally misidentifies fine-grained attributes such as color.

89 One example (black-footed albatross) illustrates a limitation: highly relevant regions sometimes correspond to background patches, leading to less informative labels. In another case (Acadian Flycatcher), the model misclassifies the species but still provides interpretable relevance maps and semantically useful labels. These findings suggest that the pipeline can enhance model interpretability, even in failure cases, and may support model refinement.

94 Quantitative Evaluation

95 We evaluate whether (1) generated crops overlap annotated body parts and (2) CLIP labels match the
96 correct part. Part annotations were mapped to CUB200 attributes (e.g., bill→beak).

Table 1: Evaluation of part matching performance.

Metric	LRP-CLIP			Random Box		
	TOP1	TOP3	TOP5	TOP1	TOP3	TOP5
Body part present in crop (%)	93.33	93.33	93.33	35.61	35.61	35.61
At least one correct CLIP label (%)	49.15	62.36	68.47	9.58	14.00	16.71

97 Table 1 shows that 93% of crops contain annotated parts and nearly 70% yield a correct CLIP label
98 (top-5). Both metrics outperform a random baseline, supporting the pipeline’s ability to generate
99 semantically valid part-level explanations.

Table 2: Evaluation of part label heuristics.

Metric	LRP-CLIP			Random Box		
	Eye	Forehead	Breast	Back	Leg	Breast
Top three matched labels in all parts (%)	23.23	21.76	11.47	18.92	16.9	15
Least three matched labels in all parts(%)	Crown 3.58	Belly 1.27	Tail 0.4	Tail 1.94	Belly 1.75	Crown 1.64

100 Table 2 shows the top three and least three distribution of the matched labels in all parts of all correctly
101 labeled crops. It shows a clear bias towards the head by LRP-CLIP, while random box focuses on
102 general big parts of the anatomy of a bird. The least three matched labels are the same between
103 LRP-CLIP and random box. This either indicates a lack of understanding for CLIP for these specific
104 labels or the little prevalence of these parts in bird pictures.

105 Discussion

106 **Explanation Module Components** LRP highlights regions that often correspond to bird parts, but
107 cropping involves trade-offs: small crops may lose context, while large ones dilute relevance. Future
108 work may explore rules that also capture negative relevance.

109 **Labeling with CLIP** CLIP labeling depends on crop quality and label sets. Using CUB-200 attributes
110 limits coverage (e.g., background context is excluded), and small or ambiguous crops often lead to
111 guesses. Nonetheless, CLIP maps abstract regions to interpretable part labels without retraining,
112 making it a lightweight semantic layer. Future work could expand label sets or preserve more image
113 context for robustness.

114 **Efficiency** The pipeline adds little overhead: one forward pass, one backward pass, and CLIP
115 embedding of crops, with labels pre-embedded for efficiency.

116 5 Future Work

117 This paper presents only first results toward enhancing LRP explanations with natural language
118 semantics. The current study is limited to a single dataset (CUB-200), one explanation method (LRP),
119 and a basic region extraction strategy. These choices allowed us to demonstrate feasibility with
120 minimal additional training effort, but they restrict generality.

121 Future work will extend the pipeline in several directions: adopting more sophisticated cropping
122 algorithms, comparing multiple explanation methods, and evaluating on additional datasets. In
123 the longer term, we aim to adapt the approach outside of vision domains like medical time-series
124 classification, where interpretability is particularly challenging yet essential for trust and adoption.

125 **References**

- 126 Reduan Achtibat, Sayed Mohammad Vakilzadeh Hatefi, Maximilian Dreyer, Aakriti Jain, Thomas
127 Wiegand, Sebastian Lapuschkin, and Wojciech Samek. AttnLRP: Attention-aware layer-wise
128 relevance propagation for transformers. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller,
129 Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp, editors, *Proceedings of
130 the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine
131 Learning Research*, pages 135–168. PMLR, 21–27 Jul 2024.
- 132 Christopher J. Anders, David Neumann, Wojciech Samek, Klaus-Robert Müller, and Sebastian
133 Lapuschkin. Software for dataset-wide xai: From local explanations to global insights with Zennit,
134 CoRelAy, and ViRelAy. *CoRR*, abs/2106.13200, 2021.
- 135 Clara Bove, Jonathan Aigrain, Marie-Jeanne Lesot, Charles Tijus, and Marcin Detyniecki. Context-
136 tualization and exploration of local feature importance explanations to improve understanding
137 and satisfaction of non-expert users. In *Proceedings of the 27th International Conference on
138 Intelligent User Interfaces*, IUI '22, page 807–819, New York, NY, USA, 2022. Association
139 for Computing Machinery. ISBN 9781450391443. doi: 10.1145/3490099.3511139. URL
140 <https://doi.org/10.1145/3490099.3511139>.
- 141 Vaibhav Dhore, Achintya Bhat, Viraj Nerlekar, Kashyap Chavhan, and Aniket Umare. Enhancing
142 explainable ai: A hybrid approach combining gradcam and lrp for cnn interpretability, 2024. URL
143 <https://arxiv.org/abs/2405.12175>.
- 144 Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori,
145 Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali
146 Farhadi, and Ludwig Schmidt. Openclip, July 2021. URL [https://doi.org/10.5281/zenodo.
147 5143773](https://doi.org/10.5281/zenodo.5143773). If you use this software, please cite it as below.
- 148 Alec Radford, Jong Wook Kim, Chris Hallacy, A. Ramesh, Gabriel Goh, Sandhini Agarwal, Girish
149 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.
150 Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- 151 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should I trust you?": Explaining the
152 predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference
153 on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, pages
154 1135–1144, 2016.
- 155 Fawaz Sammani and Nikos Deligiannis. Interpreting and analysing clip’s zero-shot image classifica-
156 tion via mutual knowledge, 2024. URL <https://arxiv.org/abs/2410.13016>.
- 157 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Mehdi
158 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski,
159 Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia
160 Jitsev. LAION-5b: An open large-scale dataset for training next generation image-text models.
161 In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks
162 Track*, 2022. URL <https://openreview.net/forum?id=M3Y74vmsMcY>.
- 163 Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi
164 Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-
165 based localization. *International Journal of Computer Vision*, 128(2):336–359, October 2019.
166 ISSN 1573-1405. doi: 10.1007/s11263-019-01228-7. URL [http://dx.doi.org/10.1007/
167 s11263-019-01228-7](http://dx.doi.org/10.1007/s11263-019-01228-7).
- 168 Shilin Sun, Wenbin An, Feng Tian, Fang Nan, Qidong Liu, Jun Liu, Nazaraf Shah, and Ping Chen.
169 A review of multimodal explainable artificial intelligence: Past, present and future, 2024. URL
170 <https://arxiv.org/abs/2412.14056>.
- 171 Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd
172 birds-200-2011 dataset. Juli 2011.