

Appendix

A. Original Rocchio formulation

The Rocchio algorithm [19] aims to generate the refined query vector $z'_q \in \mathbb{R}^d$, by fusing representations of original query z_q with relevant and irrelevant (or positive and negative) feature vectors z_p and z_n obtained in the current stage of retrieval. This updated representation is then used to perform another round of retrieval. The Rocchio query refinement rule can formally be defined as follows:

$$z'_q = \alpha z_q + \beta z_p - \gamma z_n \quad (9)$$

The positive vector is computed from the representations of the most relevant candidates, whereas the negative features are aggregated from the irrelevant ones. Specifically, the aggregated positive and negative representations are computed by averaging:

$$z_p = \frac{1}{|C_r|} \sum_{i \in C_r} z_i \quad (10)$$

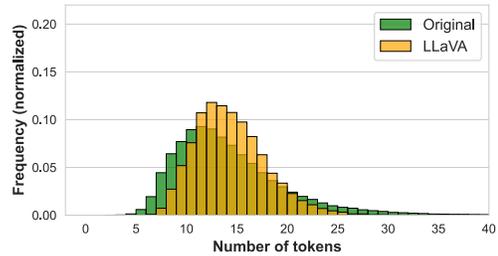
$$z_n = \frac{1}{|C_{nr}|} \sum_{i \in C_{nr}} z_i \quad (11)$$

where C_r and C_{nr} are the sets of the most relevant and non-relevant candidates with cardinalities $|C_r|$ and $|C_{nr}|$. In practice, the algorithm is often applied with $K = |C_r| = |C_{nr}|$, i.e., using the top-K candidates with the highest and lowest relevance with respect to the query.

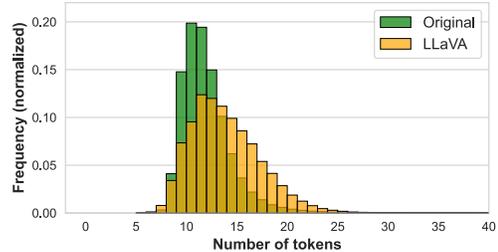
B. Generating Synthetic Captions with LLaVA-1.5

We utilized LLaVA-1.5-7B⁶ [14], quantized in 8-bit, to generate a synthetic caption for each image in both Flickr30K and COCO. Given that these captions would later be used to obtain CLIP and BLIP-2 representations – which typically rely on short and descriptive textual inputs [20, 31] – we adopted a two-step prompting strategy to balance diversity and conciseness. First, we randomly selected a prompt from a predefined pool of five sentence-level description prompts, ensuring variability in the generated responses. Second, we provided a follow-up prompt explicitly instructing the model to generate a single-sentence description. The prompts and conversation template are shown in Table 7. Figure 6 shows the distributions of the lengths of the original and generated captions. As observed, the length distribution of generated captions differs slightly from that of human-written ones while still avoiding extreme values and remaining within a reasonable range.

⁶We used open-source checkpoints available via HuggingFace: <https://huggingface.co/llava-hf/llava-1.5-7b-hf>

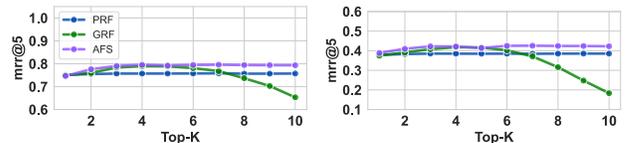


(a) Flickr30K



(b) COCO

Figure 6. Distribution of original and LLaVA-generated caption lengths.



(a) Flickr30K

(b) COCO

Figure 7. **Retrieval with Rocchio: top-K items.** Retrieval metrics with a varying number of items in relevance feedback.

Dataset	Feedback	0.05	0.10	0.25	0.50
Flickr	GRF	0.716	0.714	0.707	0.709
	PRF	0.669	0.670	0.669	0.671
COCO	GRF	0.330	0.331	0.327	0.327
	PRF	0.295	0.295	0.295	0.294

Table 6. Retrieval performance (Hits@1) with different values of temperature.

C. Rocchio Temperature

In Equation 1, we define positive and negative weights for each item in relevance feedback. Specifically, we scale similarities between each item and the corresponding query representation by temperature τ to compute the weight of each item. These weights are then used in Equation 2 to guide query representations. Table 6 shows that using lower temperature values, which sharpen the weight distribution, leads to better retrieval performance.

Step	Role	Prompt
0		<p><prompt1> is randomly selected from the pool:</p> <ul style="list-style-type: none"> • "Describe the image and main visual features in one sentence." • "Generate a short caption for this image, focusing on the visual details." • "Provide a concise description, one sentence, of the image visual features and surroundings." • "Write a brief caption, one sentence, that describes the image visual features." • "Summarize the image visual features in a precise caption, maximum one sentence." <p><prompt2>:</p> <ul style="list-style-type: none"> • "Pay attention to the visual settings and details on the image. Write exactly one sentence under 10 words."
1	User Assistant	<p>USER: <image><prompt1> ASSISTANT: USER: <image><prompt1> ASSISTANT: <answer1></p>
2	User Assistant	<p>USER: <image><prompt1> ASSISTANT: <answer1>USER: <prompt2> ASSISTANT: USER: <image><prompt1> ASSISTANT: <answer1>USER: <prompt2> ASSISTANT: <answer2></p>

Table 7. Conversation template for two-step prompting; <answer2> is the generated synthetic caption for input image <image>.

D. Number of Candidates for Relevance Feedback

We evaluate the impact of the number of items in relevance feedback. According to Figure 7, the optimal performance in MRR@5 for both GRF and AFS is achieved with 4-5 relevant items on both datasets. Further increase of the relevance set leads to decay in MRR for GRF.

E. Attentive Feedback Summarizer

E.1. Attention Blocks Architecture

Section 3.3 presents an overview of the proposed attentive feedback summarizer architecture. This section provides a detailed description of how image patches and caption tokens are processed by the summarizer.

The first component of the AFS architecture is a cross-attention module (Figure 8) designed to fuse query representations with relevance feedback. Specifically, cross-attention queries are obtained for the token features of the user query. We denote these token-level representations of input query q with s tokens as $\mathbf{h}_q = [\mathbf{h}_{q,\text{CLS}}, \mathbf{h}_{q,1}, \mathbf{h}_{q,2}, \dots, \mathbf{h}_{q,j}, \dots, \mathbf{h}_{q,s}] \in \mathbb{R}^{s \times d_t}$, where d_t is the token dimensionality and $\mathbf{h}_{q,\text{CLS}}$ is a learnable CLS token. Furthermore, keys and values are obtained from relevance feedback, comprising patch- and word-level representations of the most relevant images and their corresponding AI-generated captions. Relevance feedback vector $\mathbf{r}_q = \text{concat}([\mathbf{r}_q^{\text{img}}, \mathbf{r}_q^{\text{cap}}])$ is a concatenated sequence of the patch and token embeddings for the top-K relevant

images and their synthetic captions projected to \mathbb{R}^{d_t} . In detail, visual relevance sequence can be defined as $\mathbf{r}_q^{\text{img}} = [\mathbf{r}_{q,1}^{\text{img-1}}, \mathbf{r}_{q,2}^{\text{img-1}}, \dots, \mathbf{r}_{q,p}^{\text{img-1}}, \dots, \mathbf{r}_{q,1}^{\text{img-K}}, \mathbf{r}_{q,2}^{\text{img-K}}, \dots, \mathbf{r}_{q,p}^{\text{img-K}}]$ for images divided into p patches $\mathbf{r}_{q,j}^{\text{img-c}} \in \mathbb{R}^{d_t}$. Similarly, the word-level relevance sequence from corresponding AI-generated captions is $\mathbf{r}_q^{\text{cap}} = [\mathbf{r}_{q,1}^{\text{cap-1}}, \mathbf{r}_{q,2}^{\text{cap-1}}, \dots, \mathbf{r}_{q,s}^{\text{cap-1}}, \dots, \mathbf{r}_{q,1}^{\text{cap-K}}, \mathbf{r}_{q,2}^{\text{cap-K}}, \dots, \mathbf{r}_{q,s}^{\text{cap-K}}]$ with captions padded to length s and $\mathbf{r}_{q,j}^{\text{cap-c}} \in \mathbb{R}^{d_t}$.

The second component of the summarizer is a standard self-attention block that communicates information between query sequence tokens. Therefore, the output sequence of these two blocks can be defined as $\tilde{\mathbf{h}}_q = \text{SelfAttn}(\text{CrossAttn}(\mathbf{h}_q, \mathbf{r}_q)) \in \mathbb{R}^{s \times d_t}$. We use CLS token $\tilde{\mathbf{h}}_{q,\text{CLS}}$ as the output of the feedback summarizer and apply linear projection to obtain $\mathbf{z}_q^{\text{CLS}} \in \mathbb{R}^d$.

E.2. AFS Representations

In this section, we visualize representations $\mathbf{z}_q^{\text{CLS}}$ learned by the AFS model. We compare two AFS variants trained with different objectives: one using only the image-based loss (l_q^{img}), and one using a combined loss from both image and caption supervision ($l_q^{\text{img}} + l_q^{\text{cap}}$), corresponding to rows 2 and 3 in Table 2. As shown in Section 5.1, the variant trained with image-based loss achieves better retrieval performance. Figure 9 shows PCA projections of the query, image and AFS embeddings in two-dimensional space for Flickr30K with CLIP-ViT-B/32. The embeddings from the combined-loss model are positioned between the text and image embedding clusters, while the image-only

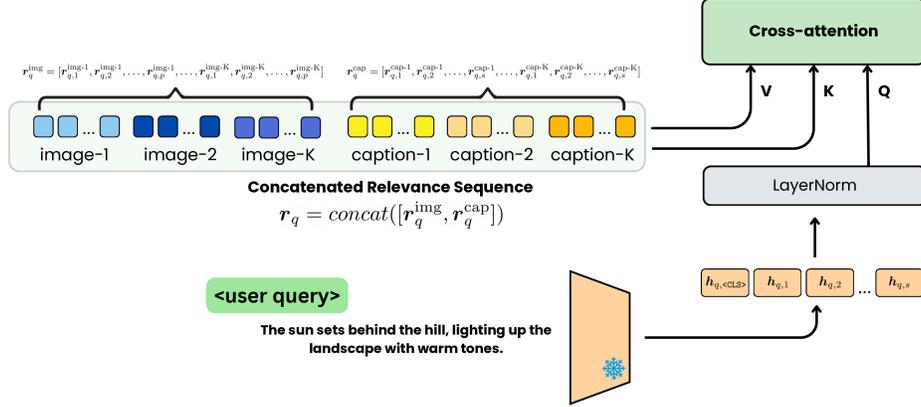


Figure 8. Cross-attention architecture.

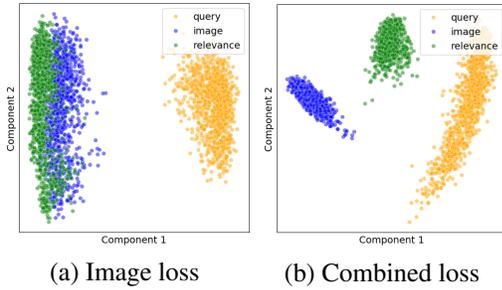


Figure 9. Images, queries and AFS representations projected with PCA.

variant produces representations that remain closer to the image features. This pattern suggests that feedback representations more aligned with image features may be more effective for refining query representations.

E.3. Cross-attention Score Aggregation

In Section 5.3 and Appendix E.5, we visualize cross-attention scores to demonstrate feedback summarization from fine-grained feedback.

In this section, we provide technical details on cross-attention score aggregation. Cross-attention scores are defined as $\mathbf{A}_q \in \mathbb{R}^{n_h \times s_q \times s_r}$, where n_h is the number of attention heads, s_q is the number of tokens in a user query, and s_r is the number of patch and token embeddings in a relevance sequence. First, we accumulate attention scores over attention heads and query tokens obtaining $\hat{\mathbf{A}}_q = [\hat{A}_{q,1}, \hat{A}_{q,2}, \dots, \hat{A}_{q,s_r}] \in \mathbb{R}^{s_r}$, where each item is a scalar corresponding to a certain patch or token in relevance sequence r_q . Then, we split $\hat{\mathbf{A}}_q$ into image and caption parts $\hat{\mathbf{A}}_q^{\text{img}} = [\hat{A}_{q,1}^{\text{img}}, \dots, \hat{A}_{q,p}^{\text{img}}, \dots, \hat{A}_{q,1}^{\text{img}}, \dots, \hat{A}_{q,p}^{\text{img}}]$ and $\hat{\mathbf{A}}_q^{\text{cap}} = [\hat{A}_{q,1}^{\text{cap}}, \dots, \hat{A}_{q,s}^{\text{cap}}, \dots, \hat{A}_{q,1}^{\text{cap}}, \dots, \hat{A}_{q,s}^{\text{cap}}]$ for K items used as relevance feedback. Finally, both sequences $\hat{\mathbf{A}}_q^{\text{img}}$ and $\hat{\mathbf{A}}_q^{\text{cap}}$ are independently normalized be-

tween 0 and 1 using min-max scaling. As a result, the normalized values are used as saliency scores for patches on relevant images and words in synthetic captions.

E.4. AFS without Synthetic Captions

As shown in Table 4, PRF does not improve over baseline retrieval models without relevance feedback (Table 4: rows 2, 7, and 12). In Sections 3.3 and 4.5, we introduced AFS as a strategy that combines PRF and GRF by aggregating local information from both image patches and synthetic captions. Here, we evaluate whether AFS can operate using retrieved image embeddings without synthetic captions. We refer to this variant as AFS-PRF.⁷

Figure 10 shows that AFS-PRF performs only slightly worse than the full AFS model while still outperforming PRF with Rocchio. This result suggests that AFS can serve as an effective tool for enabling pseudo-relevance feedback.

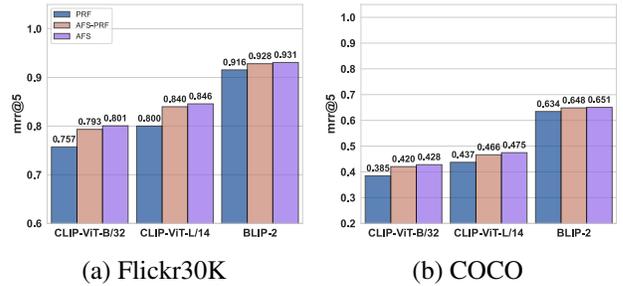


Figure 10. AFS-PRF Evaluation. Comparison between AFS without generated captions (AFS-PRF), its full version, and Rocchio-based PRF.

⁷We did not re-train the models using images only; instead, we used the same trained models from previous sections while ignoring synthetic captions at inference.

E.5. Cross-attention Visualization Examples

We provide more cross-attention visualizations obtained from the AFS model, supplementing Section 5.3. In Figures 14 and 15, saliency maps are obtained for CLIP-ViT-B/32 and CLIP-ViT-L/14, respectively. In both cases, the ground-truth images were not among the top-5 initially retrieved items and, hence, were excluded from relevance feedback aggregation. Nevertheless, the obtained attention maps highlight semantically meaningful regions, and the refined query representations lead to better retrieval performance. Furthermore, Figures 16 and 17, both using CLIP-ViT-L/14, show scenarios where the ground-truth image ranks improved after applying AFS.

F. Multi-turn Retrieval Results

This section presents the extended results of the multi-turn retrieval with relevance feedback aggregation. Figures 11 and 12 show the performance metrics achieved across retrieval rounds for the Flickr30K and COCO datasets. The obtained results complement the findings discussed in Section 5.2. Specifically, explicit feedback with ground truth captions continuously improves retrieval performance according to all metrics. Generative relevance feedback, however, leads to performance degradation starting from the third retrieval round. The attentive feedback summarizer demonstrates an increase in performance at round 2, gradually converging in rounds 3-5, avoiding query drift without ground-truth captions.

G. Interactive Retrieval Demo

In Section 6, we outlined future directions for combining relevance feedback techniques with user interactions to create more interactive retrieval systems. While some recent work explores chat-based interactions via LLMs [8, 32], we investigate an alternative approach based on direct visual interaction with images. Specifically, we developed a prototype interface where users can retrieve images based on a textual query. The interface, then, allows users to annotate retrieved image regions by drawing bounding boxes indicating relevance or irrelevance to their search intent (Figure 13).

We explored how the bounding boxes could be integrated into the inference process of the AFS model working without synthetic captions (Appendix E.4). Specifically, we modified the cross-attention scores based on a simple heuristic: attention weights were increased for image patches corresponding to regions marked as relevant by the user and decreased for those identified as irrelevant. The magnitude of the increase is a hyperparameter and can be defined through configurations. These adjustments were applied at inference time without requiring any additional training or fine-tuning of the model.

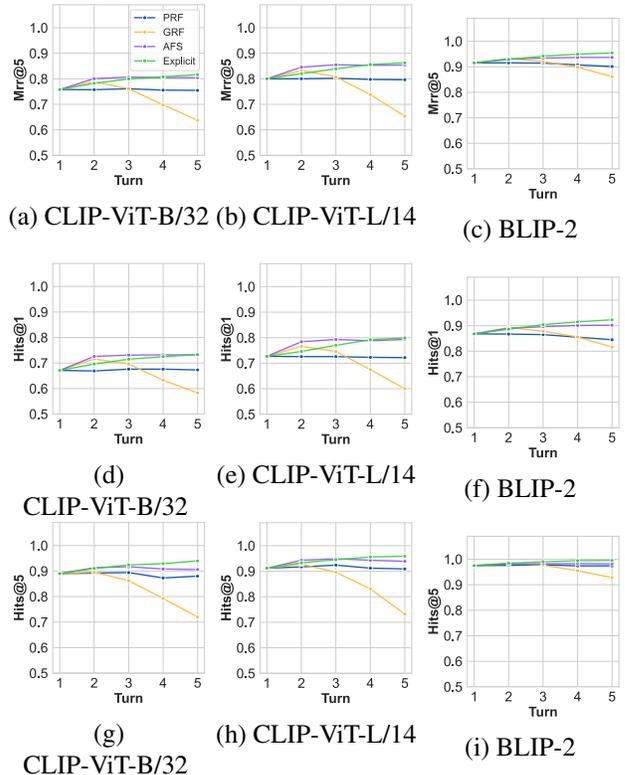


Figure 11. Multi-turn retrieval performance with relevance feedback on Flickr30K.

As shown in Figure 13, this approach effectively shifts the model’s attention toward user-indicated regions, enabling more targeted retrieval responses. While our method is based on fixed weight adjustments, it demonstrates how user input can guide attention in relevance-based models. In the future, more sophisticated techniques could be explored. For instance, it is possible to add learnable token-level embeddings based on user feedback. This would require training the model with region-level relevance annotations, which are currently not available.

The code for the prototype is available at: https://github.com/bulatkh/visualref/tree/wacv_demo

H. Limitations

Our study systematically evaluated relevance feedback strategies with pre-trained VLM backbones. However, we do not compare against alternative query adjustment methods, such as query rewriting or prompt engineering with LLMs-in-the-loop. Future work can address this gap and explore the interplay between representation-level refinement and natural language query adjustments. This direction is particularly relevant for dialogue-aware image retrieval, where systems must model contextual coherence

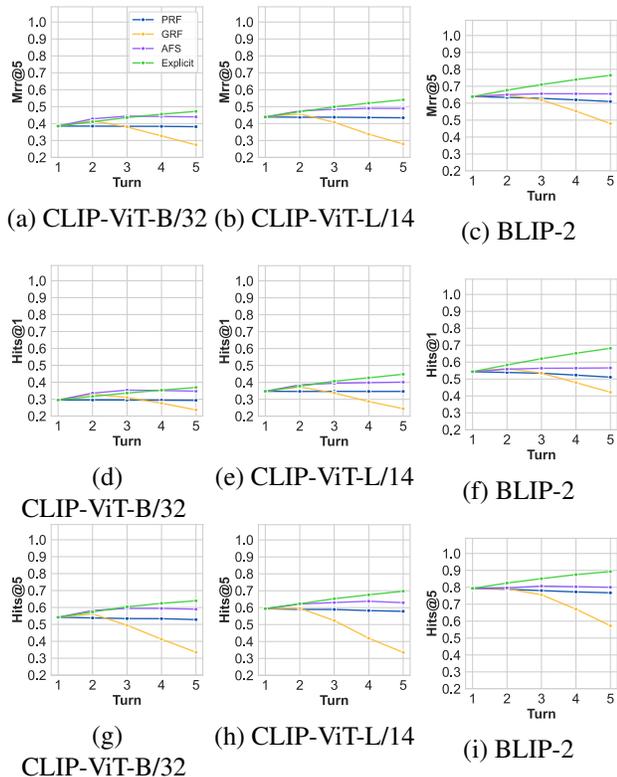


Figure 12. Multi-turn retrieval performance with relevance feedback on COCO.

across multiple modalities and turns, resolve ambiguous references, and encode rich dialogue context for effective prompting of VLMs, which often struggle with non-descriptive queries [20].

Further, GRF and AFS strategies introduce some additional overhead. GRF requires running an image captioning model in the background to generate captions for candidate images. AFS, on the other hand, operates at runtime, although its size remains modest (under 20 million parameters), especially compared to LLMs, which can contain billions of parameters. Therefore, we suggest that developers carefully weigh the trade-off between the retrieval gains offered by relevance feedback and the computational (and environmental) costs associated with deploying these methods.

Finally, we tested the proposed relevance feedback strategies on general-purpose image retrieval datasets. In future work, we are planning to evaluate these methods in domain-specific retrieval tasks.

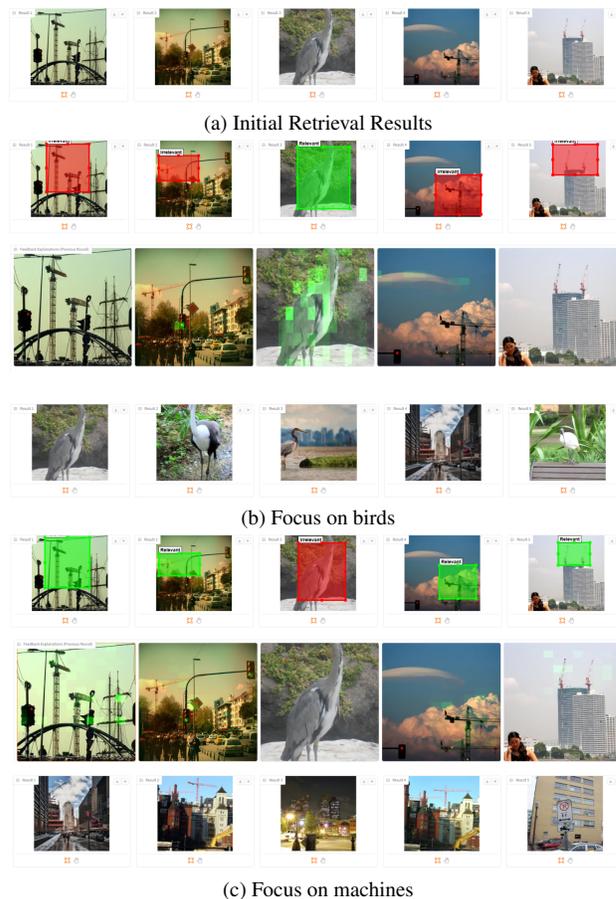


Figure 13. **Interactive Retrieval Demo.** (a) Initial retrieval results for the textual query "Crane in the city". (b) The user specifies their intent by marking the bird as relevant and industrial cranes as irrelevant. (c) The user instead focuses on industrial cranes. In both (b) and (c), the first row shows examples of manual relevance annotations, the second row visualizes cross-attention scores after applying the feedback, and the third row presents the updated retrieval results using AFS with user signals.

Query:

players in the red and white uniforms and bystanders are standing on a soccer field

Top 5 retrieved images with CLIP-ViT-B/32



Cross-attention visualization

Patch-level scores



Token-level scores

team	is	sitting	on	the	field	with	some	of	them	wearing	red	jerseys						
soccer	game	is	taking	place	on	a	field	with	players	in	uniforms	and	a	soccer	ball			
group	of	soccer	players	wearing	yellow	jerseys	are	standing	on	a	field							
soccer	player	in	a	yellow	shirt	and	green	shorts	is	standing	on	a	field	in	front	of	a	goal
a	large	crowd	of	people	walking	down	a	street										

Top 5 retrieved images after feedback aggregation



Figure 14. Cross-attention visualization with CLIP-ViT-B/32. This example shows a case when a ground-truth image was not used for relevance feedback aggregation, i.e., it was not initially among the top 5 retrieved images.

Query:

male in the street wearing a blue plaid shirt waiting for the light to turn

Top 5 retrieved images with CLIP-ViT-L/14



Cross-attention visualization

Patch-level scores



Token-level scores

man	is	in	a	red	shirt	is	sitting	on	the	sidewalk	looking	at	his	cell	phone			
man	is	a	man	walking	is	standing	on	a	sidewalk	next	to	a	car					
man	is	a	man	wearing	a	plaid	shirt	is	walking	in	front	of	a	blue	door			
man	wearing	a	plaid	shirt	is	walking	a	street										
man	is	a	man	is	looking	on	his	cell	phone									

Top 5 retrieved images after feedback aggregation



Figure 15. Cross-attention visualization with CLIP-ViT-L/14. The ground-truth image was not used for relevance feedback aggregation, i.e., it was not initially among the top 5 retrieved images.

Query:

a man dressed in a checkered shirt and black pants performs an aerial stunt on a bicycle

Top 5 retrieved images with CLIP-ViT-L/14



Cross-attention visualization

Patch-level scores



Token-level scores

men	is	doing	a	stunt	on	a	sidewalk												
man	is	performing	a	trick	on	a	bicycle	with	his	feet	on	the	wheel						
man	on	a	motorcycle	is	performing	a	trick	in	front	of	a	crowd							
man	is	performing	a	trick	on	a	bicycle	flying	the	air	while	riding	the	bike					
man	is	riding	a	bike	in	the	air	with	a	blue	sky	and	clouds	in	the	background			

Top 5 retrieved images after feedback aggregation



Figure 16. Cross-attention visualization with CLIP-ViT-L/14. AFS increases the rank of the ground-truth image.

Query:

two horse racing jockeys, one in checkered blue and red and the other in orange and brown, are racing against a blurry background

Top 5 retrieved images with CLIP-ViT-L/14



Cross-attention visualization

Patch-level scores



Token-level scores

cowboy	riding	a	horse	in	a	dirt	field													
man	riding	a	brown	horse	is	racing	against	other	horses	in	a	dirt	track							
race	car	with	the	number	2	6	on	it	is	driving	on	a	track							
greyhound	dogs	wearing	red	and	green	vests	are	running	a	dirt	track									
jockey	riding	a	horse	with	a	red	and	white	blanket	on	it									

Top 5 retrieved images after feedback aggregation



Figure 17. Cross-attention visualization with CLIP-ViT-L/14. AFS increases the rank of the ground-truth image.