

Imagine and Seek: Improving Composed Image Retrieval with an Imagined Proxy

You Li Fan Ma Yi Yang[†]
ReLER, CCAI, Zhejiang University, Zhejiang, China

[†] Corresponding author

{uli2000, mafan, yangyics}@zju.edu.cn

Abstract

The Zero-shot Composed Image Retrieval (ZSCIR) requires retrieving images that match the query image and the relative captions. Current methods focus on projecting the query image into the text feature space, subsequently combining them with features of query texts for retrieval. However, retrieving images only with the text features cannot guarantee detailed alignment due to the natural gap between images and text. In this paper, we introduce **Imagined Proxy for CIR (IP-CIR)**, a training-free method that creates a proxy image aligned with the query image and text description, enhancing query representation in the retrieval process. We first leverage the large language model’s generalization capability to generate an image layout, and then apply both the query text and image for conditional generation. The robust query features are enhanced by merging the proxy image, query image, and text semantic perturbation. Our newly proposed balancing metric integrates text-based and proxy retrieval similarities, allowing for more accurate retrieval of the target image while incorporating image-side information into the process. Experiments on three public datasets demonstrate that our method significantly improves retrieval performances. We achieve state-of-the-art (SOTA) results on the CIR dataset with a Recall@K of 70.07 at K=10. Additionally, we achieved an improvement in Recall@10 on the FashionIQ dataset, rising from 45.11 to 45.74, and improved the baseline performance in CIRCO with a mAPK@10 score, increasing from 32.24 to 34.26.

1. Introduction

Imagination is the wellspring of human creativity. When you see a picture of an adorable cat, you might envision it playing with other cats, wearing cool sunglasses and a cape, or floating in a space capsule. This ability to imagine and pursue beautiful visions drives continuous advancement in fields like image generation [3, 16, 19, 25, 29, 32, 35, 46,

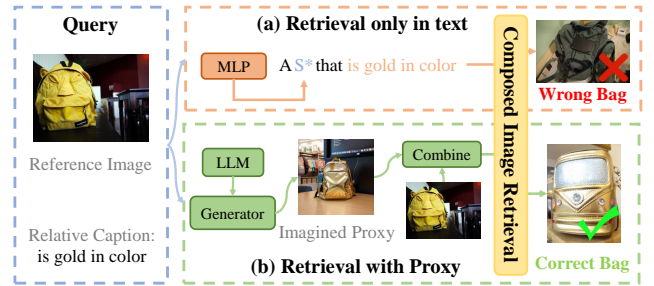


Figure 1. **Comparison of composed image retrieval between text-only retrieval and our methods.** Traditional methods perform retrieval only in the text space, where CLIP text features may overlook some important semantic information. In contrast, our approach generates imagined proxy features, providing additional information that is often overlooked in text-only retrieval, thereby improving retrieval accuracy.

49] and image retrieval [1, 6, 7, 17, 24, 27, 34, 40, 43, 44].

This rich imagination gains new significance in *Composed Image Retrieval (CIR)* [27], where users express their imagination of the query image through text, searching for content that aligns with their preferences. Specifically, the CIR task requires users to retrieve target images from a database that match both the content of a query image and relative caption. With text modality to flexibly control retrieval targets, *CIR enables the search process akin to the user’s specific imagination after viewing the query image*. This **multi-modal, imagination-like** retrieval task greatly expands the applications of image retrieval, with wide-ranging uses in fields such as visual search [28], object localization [18, 21], and re-identification [38, 45].

Traditional solutions for CIR [4, 5, 10–15] are constrained by the scarcity of training data (e.g. 46.6k for FashionIQ [40] dataset and 28.8k for CIR dataset [27]), as collecting suitable triplets (query image x_q , relevant text x_c , and target image x_t) from the web incurs high labor costs. Therefore, we focus on the ZS-CIR task, aiming to bypass the need for training on triplet datasets and achieve gen-

eralized CIR retrieval capabilities. As shown in Fig.1 (a), traditional ZS-CIR methods typically employ lightweight projection models to map CLIP’s [31] image features into the text features [7] for retrieval. *The transformation inevitably loses certain image information, and the coarse, easily confusable text features fall short in fine-grained, complex imagination scenes.* Although recent approaches attempt to leverage large language models (LLMs) [2, 30] to generate various descriptions of target images [20, 43, 44], they overlook the potential for direct imagination on the image side. As shown in Fig.1 (a), the text features did not effectively capture the attribute ‘golden’, resulting in sub-optimal retrieval results.

The rapid advancement of controllable generative models [8, 9, 23, 33, 37, 41, 42, 48] has empowered users with the freedom to create images that align with their imagined visions. Generating images that match both the content of query images and relative captions aligns closely with the goals of composed retrieval tasks, and the imagined images (we called **Retrieval Proxy**) could provide additional information—such as style, instance attributes, and spatial relationships—that is often overlooked in text-based retrieval with complex captions. This raises an intriguing question: *How could we harness the imaginative power of generative models to improve retrieval performance?*

We thus propose IP-CIR (Imagined Proxy for CIR), a training-free, plug-and-play method that harnesses the power of imagination for any retrieval method. Firstly, by utilizing LLMs to understand both the image content and the relative captions, we generate a suitable image layout. Using controllable generation methods, we create high-quality, fine-grained proxy images. Next, we compose the proxy images, query images as well as semantic perturbation into robust proxy features that are more suitable for retrieval. Finally, we propose a balance metric to combine the proxy and text retrieval similarities, allowing the integration of our method with any retrieval methods.

We conduct experiments on three public datasets: the CIRR [27], CIRCO [6], and FashionIQ [40]. We combine our method with several baselines, the observed performance improvements confirm that our method supplements the retrieval process with valuable information and could enhance the performance of any other retrieval methods.

Our contributions are summarized as follows:

- 1) We propose using generative models’ imagination capabilities to create proxy images that align with the content of the query image and relative captions to improve retrieval performance.
- 2) We introduce a proxy image generation framework. Using the LLM-generated object layout and the query image, we create imagined proxies through MIGC++.
- 3) We propose integrating proxy images with the query image and semantic perturbation into more robust proxy fea-

tures and introducing a balance metric to combine the proxy and baseline similarities.

- 4) We conducted experiments on three public CIR datasets. The improvement over baseline results, as well as achieving SOTA results demonstrates the potential of constructing imagined proxies in retrieval.

2. Related work

2.1. Zero-shot Composed Image Retrieval

Early work on ZS-CIR, such as Pic2word [34], SEALRE [7], attempted to project the original image into the text space using lightweight projection models. This was often done within fixed templates (e.g., “a photo of [\$] with cond”), where the query images and relative captions were combined. LinCIR [17] further advanced this approach by defining keywords as adjectives and nouns in the text, thus overcoming the limitations of fixed templates and training more powerful projection models. In addition, some methods have embraced LLMs [2, 30]. For instance, CIReVL [20] proposed generating captions for reference images and then using LLMs to reassemble them into target texts modified by relative captions for retrieval. LDRE [43] took this further by using LLMs to construct various edited captions and combining them based on similarity, addressing the issue of fuzzy retrieval. While these methods improved retrieval accuracy by introducing images into the text space, they all overlooked the potential of imagining proxy images to provide additional information in image size.

2.2. Controllable Text-to-Image Generation

In recent years, with the development of diffusion models [19, 35], high-quality images that meet specific requirements can now be generated. Using mainstream generators such as Stable Diffusion [32], users can generate images aligned with text descriptions. Although these foundational generators offer a certain level of control, their performance in fine-grained attribute and orientation control is limited, primarily due to the attribute ambiguity in the CLIP features of the text [47]. As a result, some efforts have focused on leveraging layout information to enhance control. For example, GLIGEN [23] introduces Gated-Self-Attention to combine positional feature information, thereby improving control over positioning. InstanceDiffusion [37] further improves attribute control on this foundation. Meanwhile, MIGC [47] adopts a divide-and-conquer approach, breaking down complex instance generation into multiple single-instance generations and incorporating Enhancement Attention and Layout Attention to improve both positional and attribute control. MIGC++ [48] further integrates ELITE [39] to transform images into CLIP-based text features and introduces an additional EA layer to process information from

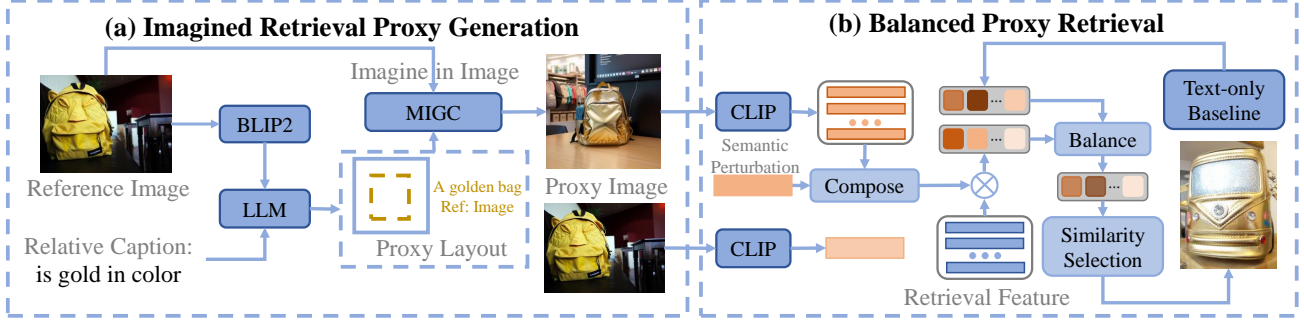


Figure 2. **Overview of our method.** (a) represents our imagined proxy generation process. We use LLM to analyze the BLIP2-generated query image captions and the relative captions and infer the proxy layout. We then use the controllable generator to imagine the proxy images. (b) represents our process of constructing a robust proxy feature, and balancing the text and proxy similarities. We integrate proxy features, query image features as well as semantic perturbations into a robust proxy feature, and propose a balance metric for retrieval.

the reference image modality, thereby incorporating image modality control in the generation process. In this paper, we aim to utilize the fine-grained control capabilities of MIGC to display certain attributes and other information that are difficult to clearly express in text at a fine-grained level within the image itself. This provides supplementary information during the retrieval process.

3. Method

In this section, we will introduce our IP-CIR, which attempts to imagine an appropriate proxy image within the image feature space, balance, and combine the retrieval results from the image side with those from the text side. Specifically, in section 3.1, we introduce the overview of our method, in section 3.2, we will describe how to construct an imagined retrieval proxy for each query. In section 3.3, we will further explain how to organize and integrate the information from the imagined retrieval proxy and present how to combine the retrieval information from both the image and text feature spaces.

3.1. Overview

The overall structure of our framework is shown in Fig.2. First, we generate the imagined proxy image for retrieval. The proxy image needs to align with the content of the query image and be as consistent as possible with the description in the relative captions. Therefore, we **first generate a reasonable object layout** and then use a controllable generation model to **imagine this layout as a concrete image**. Specifically, we use BLIP2[22] to generate a set of captions for each query. Combine with the relative captions, we instruct the LLM to reason out a suitable spatial layout of the target image, so that, with **as much accuracy as possible** in terms of *scene, instance attributes, and quantity* with relative caption, we can imagine a proxy image of the retrieval target to provide additional informa-

tion for the retrieval process. In Fig.2 (b), we compose the imagined proxy with query image and text features, forming a robust retrieval proxy. The similarity of robust proxy is further balanced with other baselines' similarity, improving the retrieval accuracy.

3.2. Imagined Retrieval Proxy Generation

Text retrieval features may struggle to precisely capture important semantics in the relevant text and the query image, especially when the caption is complex. **Generating an image that meets the retrieval requirements**, which we refer to as **imagining the retrieval proxy**, could provide more valuable details to the retrieval process.

Imagine Proxy Layout. The first step of imagining a suitable proxy relies on *envisioning the overall layout*, which requires a clear understanding of the semantic content and relationship between the query images and relative captions. We thus leverage the robust interpretive and reasoning abilities of LLMs[2, 30] to generate a target proxy layout. As illustrated in Fig.2 (a), we organize the BLIP2 generated query image captions with the relative caption into the fixed form ‘Given an image of {caption}, we show {rule}’, based on which the LLM is instructed to infer the reasonable overall scene and the layout of the proxy image. The layout includes a detailed description of each instance as well as the bounding box (bbox) coordinates. Besides, certain instances may require detailed visual information from text or image modality. For instance, when retrieving an image of a dog wearing a hat, the attributes of the dog in the proxy image come from the dog in the query image, while the attributes of the hat come from the text modality information. We thus require the LLM to determine the reference modality for each instance.

Proxy Image Generation. The next step is to *imagine the proxy* based on the generated layout. Since the CIR task is a **fuzzy retrieval task** [43], this *reduces the need to main-*

Benchmark		CIRCO				CIRR						
Metric		mAP@K				Recall@K				Recall _{Subset} @K		
Backbone	Method	k=5	k=10	k=25	k=50	k=1	k=5	k=10	k=50	k=1	k=2	k=3
ViT-L/14	SEARLE	11.68	12.73	14.33	15.12	24.24	52.48	66.29	88.84	53.76	75.01	88.19
	SEARLE-OTI	10.18	11.03	12.72	13.67	24.87	52.31	66.29	88.58	53.80	74.31	86.94
	CIReVL	18.57	19.01	20.89	21.80	24.55	52.31	64.92	86.34	59.54	79.88	89.69
	LDRE	23.35	24.03	26.44	27.50	26.53	55.57	67.54	88.50	60.43	80.31	89.90
	LDRE + IP-CIR	26.43	27.41	29.87	31.07	29.76	58.82	71.21	90.41	62.48	81.64	90.89
	(vs. LDRE)	+3.08	+3.38	+3.43	+3.57	+3.23	+3.25	+3.67	+1.91	+2.05	+1.33	+0.99
ViT-G/14	CIReVL	26.77	27.59	29.96	31.03	34.65	64.29	75.06	91.66	67.95	84.87	93.21
	LinCIR	20.34	21.85	23.98	25.14	34.75	64.12	75.93	93.42	62.36	81.78	91.28
	LinCIR + IP-CIR	25.70	26.64	29.09	30.13	35.37	64.70	76.15	93.71	62.58	81.74	91.35
	(vs. LinCIR)	+5.40	+4.79	+5.11	+4.99	+0.62	+0.58	+0.22	+0.29	+0.22	-0.04	+0.07
	LDRE	31.12	32.24	34.95	36.03	36.15	66.39	77.25	93.95	68.82	85.66	93.76
	LDRE + IP-CIR	32.75	34.26	36.86	38.03	39.25	70.07	80.00	94.89	69.95	86.87	94.22
	(vs. LDRE)	+1.63	+2.02	+1.91	+2.00	+3.10	+3.68	+2.75	+0.94	+1.13	+1.21	+0.46

Table 1. Quantitative results in CIRCO and CIRR datasets.

tain the query image’s instance ID during the generation process. So we follow the idea of MIGC++ [48], incorporating the query image features transformed by ELITE [39] into MIGC, thereby introducing some degree of the query image information into the Proxy image. Besides, an instance can be influenced by both the reference image and the relative caption. For example, if a user envisions a hat in the query image with a specific pattern, the hat should retain its original shape but include the target pattern. So, we copy the layout of the instance with reference image modality and render it based on both image and text, ensuring the proxy image reflects the text’s modifications.

3.3. Balanced Proxy Retrieval

Once the proxy image is generated, the key issue lies in how to leverage this proxy image in the retrieval. A straightforward method is to compute the proxy similarity of the retrieval database, weighted sum with the text similarity. However, as shown in the second line of Fig.6, directly using the CLIP features of the proxy image (+PI) can lead to an excessive focus on certain elements, like the beach background. This may cause the results to overly emphasize the visual information in the image while neglecting some details from the text. So in this section, we propose to construct a robust retrieval proxy and balance the retrieval result from image and text.

Constructing Robust Proxy. Inspired by LDRE [43], we use the LLM [2, 30] to infer a set of target captions f_t . With the BLIP2 [22] generated query captions f_o , we are able to derive features that include semantic perturbation $f_s = f_t - f_o$, which represents a reasonable perturbation of the proxy features based on the edit direction. We be-

lieve that **introducing semantic perturbation** can help us to some extent **ignore overly emphasized information in the proxy images**. Besides, the proxy image may *lose some of the query image information*, so we also incorporate the features of the query image to compensate for this loss of semantics. As shown in Fig.2 (b), we combine the query image’s features f_q , the proxy image’s feature f_p with the semantic perturbation f_s into the robust proxy f_{RP} :

$$f_{RP} = f_p + \frac{\max(f_p)}{\max(f_q)} f_q + \frac{\max(f_p)}{\max(f_s)} * f_s. \quad (1)$$

Additionally, the weights of each component can be adjusted to control retrieval bias based on the quality of the proxy image and specific retrieval goals. For more details, refer to the supplementary materials.

Balancing Retrieval Results. After obtaining the retrieval feature in text and image spaces, a critical question arises: *How should we balance the two retrieval results?* The proxy features we obtain exist in the image feature space, making it difficult to directly perform weighted summation with the baseline text features. Therefore, we choose to combine their similarities. A simple approach is to take the average or weighted sum, however, in extreme cases, when we have a case with a 0.999 text similarity but only a 0.001 proxy similarity, and another case with both a 0.5 base similarity and a 0.5 proxy similarity, their averages are both 0.5, which *could lead to retrieving overly extreme results*. We thus propose a balanced metric for combining these results. Specifically, given the baseline similarity S_t and proxy similarity S_p , the final similarity S_f for retrieval is as follows:

Type		Shirt		Dress		Toptee		Average	
Backbone	Method	R@10	R@50	R@10	R@50	R@10	R@50	R@10	R@50
ViT-G/14	Pic2Word	33.17	50.39	25.43	47.65	35.24	57.62	31.28	51.89
	SEARLE	36.46	55.35	28.16	50.32	39.83	61.45	34.81	55.71
	LinCIR	46.76	65.11	38.08	60.88	50.48	71.09	45.11	65.69
	LinCIR + IP-CIR	48.04	66.68	39.02	61.03	50.18	71.14	45.74	66.28
	(vs. LinCIR)	+1.28	+1.57	+0.94	+0.15	-0.30	+0.05	+0.63	+0.59

Table 2. Quantitative results in FashionIQ dataset.

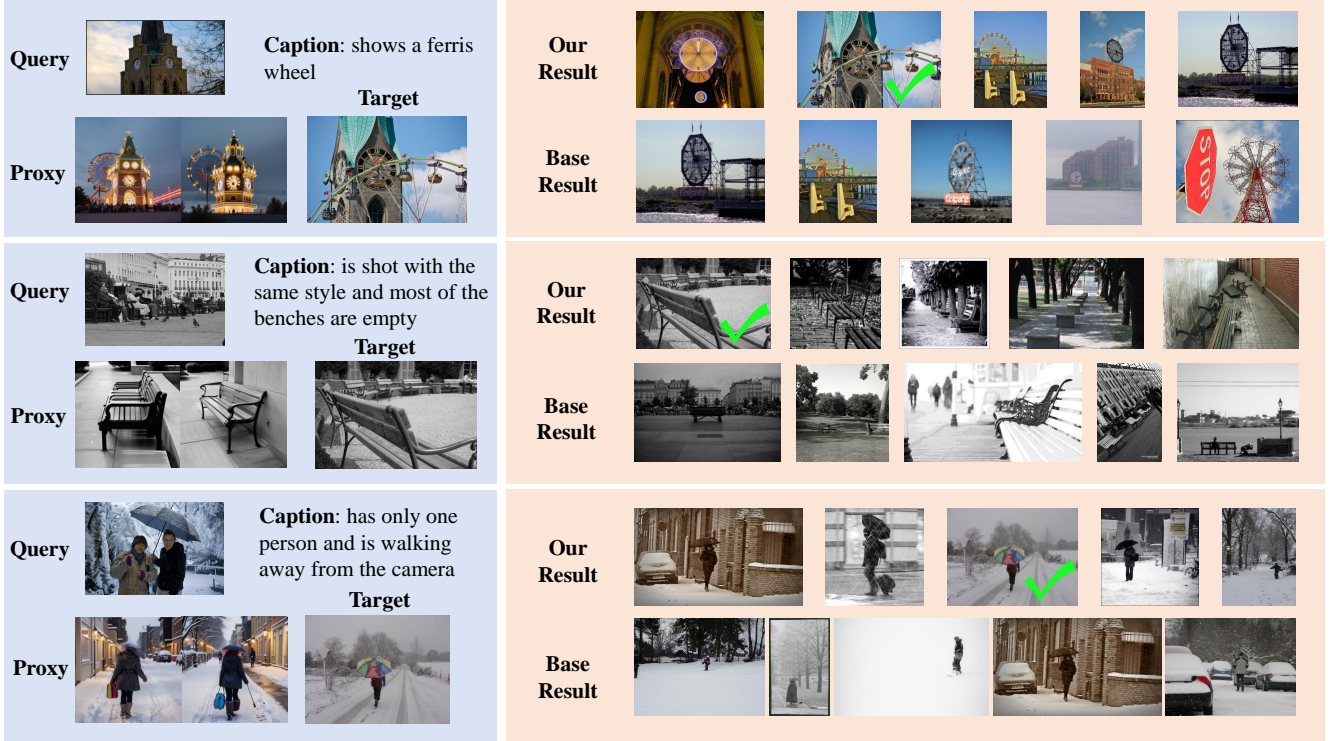


Figure 3. **Qualitative results of our method.** We conducted experiments on the CIRCO validation dataset to observe in which cases our method improves retrieval results. In the blue section on the left, we display the query information used for retrieval, the ground truth target image, and our generated proxy image features. ‘Query’ represents the input query image, ‘Caption’ represents the relative text, and in the ‘Proxy’ section, we show two generated imagined proxies. In the red section on the right, the top and bottom rows display the top-5 retrieval results enhanced by our imagined proxies and the baseline’s top-5 retrieval results, respectively.

$$S_f = \lambda S_t + (1 - \lambda) S_b, S_b = S_t * S_p, \quad (2)$$

where S_b represents the balanced similarity, λ is a balancing parameter. We believe that the balance metric enables us to focus more on cases that perform well in both similarities.

4. Experiments

In this section, we demonstrate the effectiveness of our method. In Section 4.1, We will provide a detailed description of our experimental setup, including the baselines, evaluation metrics, and implementation details. In Section 4.2, we compare our method with other state-of-the-art ap-

proaches. We explore the improvements in retrieval performance by presenting some visual examples, highlighting the impact of our method in Section 4.3, and in Section 4.4, we conduct ablation studies to investigate the contribution of each component of our method, as well as the influence of various parameter settings.

4.1. Experimental Setup

Datasets. We validate the performance of our method by combining them with several state-of-the-art methods on three datasets: CIRR [27], CIRCO [6], and FashionIQ [40]. CIRR contains 21,552 real images from the NLVR [36] dataset. CIRCO is built from the COCO 2017 unlabeled

set [26], includes a validation set with 220 queries and a test set with 800 queries. FashionIQ [40] is a dataset dedicated to the fashion domain, organized into three distinct subcategories: Dress, Shirt, and Toptee. It contains 30,135 query triplets and 77,683 images available for retrieval.

Baseline. We selected LDRE[43] and LinCIR[17] as our baselines. For each query, we generate 5 proxy images and obtain the enhanced retrieval results. We compare the improved results with some of the most recent ZS-CIR results, such as Searle[7], CIReVL[20].

Evaluation Metrics. In **CIRCO**, each query has multiple target images, so we calculate mAP@k to offer a more detailed assessment, where $k \in \{5, 10, 25, 50\}$ indicates the number of top-ranked retrieval results. For **CIRR**, we follow the original benchmarks and use Recall@k (with $k \in \{1, 5, 10, 50\}$) as the primary metric. We also assess performance in a subset setting, which we denote as Recall_{Subset}@K (with $k \in \{1, 2, 3\}$). On **FashionIQ**, we also use Recall@K (with $k \in \{10, 50\}$) for evaluation.

Implement Details. We primarily conducted experiments using two backbone models: CLIP-L and CLIP-G. We used Qwen1.5-32B[2] as the LLM in all the experiments, which yielded results comparable to those reported in the original papers. For generating proxy images, we use MIGC[47] with SD1.5 as the backbone. We use Qwen1.5-32B as our LLM to generate the layout for each query and then generate 5 proxy images based on the query image and layout. All the experiments are conducted by Pytorch with one NVIDIA A6000 GPU.

4.2. Quantitative Results

CIRCO. We present the experimental results on the CIRCO test dataset on LDRE[43] and LinCIR[17]. In the left half of Tab.1, we observe that our proxy image can improve the retrieval performance, which an approximately 5-point improvement in mAP at k=5 in LinCIR, and about 2-point improvement in LDRE. These results effectively demonstrate the validity and impact of our method.

CIRR. In the right half of Tab.1, we present the experimental results on the CIRR test dataset. Our proxy image significantly improved the LDRE and LinCIR performance. The performance improvements across different methods and backbones indicate that our approach can be applied to various text-based retrieval methods.

FashionIQ. In Tab.2, we present the results on the FashionIQ val dataset. Since LDRE has not open-sourced the code for the FashionIQ, we used LinCIR as the baseline. Although we observe a drop in R@10 in the Toptee, our method improves the retrieval accuracy in the Shirt and dress Subset, thus improving the average performance.



Figure 4. **Qualitative results of our method.** We present the result on CIRCO validation dataset with top-1 retrieval results. Baseline methods in text-based retrieval overlook certain attributes and details in complex related texts, resulting in suboptimal retrieval outcomes. In contrast, our approach can construct suitable proxies on the image side, supplementing this missing information and yielding improved retrieval results.

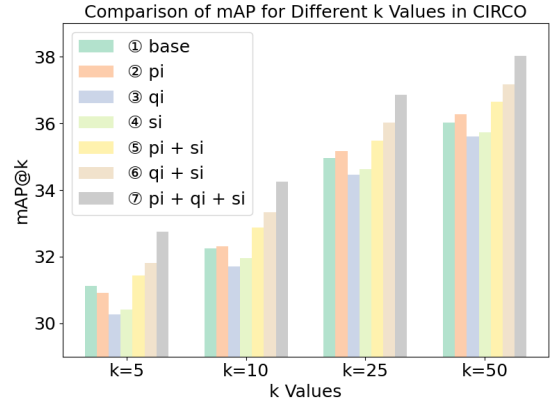


Figure 5. **Ablation results on the composition of robust proxy features** in the CIRCO dataset. *pi* indicates proxy features, *si* represents semantic perturbation, and *qi* indicates the query features.

4.3. Qualitative Results

We conducted experiments on the CIRCO validation dataset to verify the improvements our method brings to retrieval.

Experiment setting. Specifically, we used LDRE as the baseline and generated 5 imagined proxy images for each query. We present two types of results. First, based on the ground truth for each query, we show in Fig.3 that our method improves the accuracy of ground truth occurrence in top-5 retrieval results. Second, disregarding the ground truth image, we demonstrate in Fig.4 that our method re-

Benchmark				CIRR					CIRCO			
Metric				Recall@K					mAP@K			
Backbone	PI	RR	BM	k=1	k=2	k=5	k=10	k=50	k=5	k=10	k=25	k=50
LDRE-G				36.15	49.49	66.39	77.25	93.95	31.12	32.24	34.95	36.03
	✓			36.02	50.02	66.51	77.52	93.74	31.06	32.41	35.19	36.35
				-0.13	+0.53	+0.12	+0.27	-0.21	-0.06	+0.17	+0.24	+0.32
	✓	✓		38.15	52.53	68.55	79.47	94.48	31.92	33.59	36.26	37.37
				+2.00	+3.04	+2.16	+2.22	+0.53	+0.80	+1.35	+1.31	+1.34
	✓	✓	✓	39.25	52.94	70.07	80.00	94.89	32.75	34.26	36.86	38.03
				+3.10	+3.45	+3.68	+2.75	+0.94	+1.63	+2.02	+1.91	+2.00

Table 3. Ablation results in CIRR and CIRCO dataset.



Figure 6. Visualize of Ablation Result. We present the ablation result on CIRCO dataset.

trieves a more suitable image in the top-1 result.

Analysis of Experimental Results on TOP-5 Retrieval.

From the visualized results in Fig.3, we can draw the following conclusions: (1) Our generated proxy images effectively preserve the style and key elements of the original image. In the first row, our image maintains the clock tower as a significant element and adheres to the related caption by generating a Ferris wheel. In the second row, our image retains the original black-and-white style. (2) The imagined proxies complement important information that may be overlooked by text alone. For instance, in the first-row example, the clock tower—a crucial concept—is missing in the text-only retrieval result, while our imagined proxy and the original image both provide this retrieval information. In the third row, our constructed image improves retrieval accuracy by incorporating an umbrella and snow.

Analysis of Experimental Results on TOP-1 Retrieval.

Cases in Fig.4 show that by providing proxy image information, our approach achieves more relevant retrievals. **Retrieval based solely on text space may only capture partial information in complex contexts**, leading to the **omission of certain details** in the images—such as the person in the first row, the tie in the second row, and the White House in the background of the third row. Our method incorporates these details in the images, supplementing and balancing this information to yield more suitable results.

4.4. Ablation Study

Effects of the Proxy Image Features. We first conduct an ablation study on the proposed proxy image features in the CIRR dataset and CIRCO dataset. We use PI (Proxy Image) refers to the approach that directly computes retrieval similarity from proxy image features, and then performs a weighted sum with the original similarity. RP (Robust

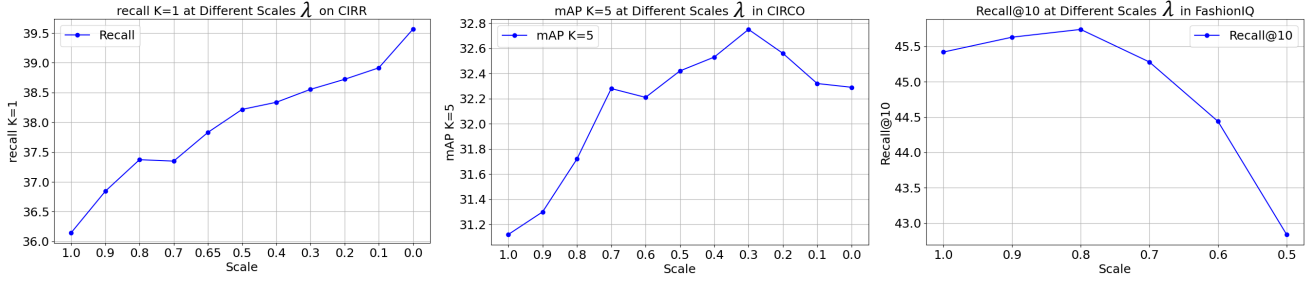


Figure 7. **The impact of weight λ in CIRCO, CIRR, and FashionIQ.** As λ decreases, the proportion of proxy information increases.

Proxy) indicates integrating the proxy image into more robust proxy features. BM (Balance Metric) refers to using the balance metric. The results from the CIRR and CIRCO are shown in Tab.3. We can draw the following conclusions: 1) The imagined proxy images may contain some noise, leading to the retrieval focus on irrelevant details. we observe a decrease in recall at $k=1$ and 50 in CIRR, and mAP $k=5$ decreases in CIRCO. 2) After integrating the proxy image into a more robust proxy feature, we see a significant improvement in both datasets, which demonstrates the effectiveness of our robust proxy features. 3) By introducing the balance metric between proxy and text similarity, we further improve the retrieval performance. This suggests that combining the information from both image and text domains leads to more accurate retrieval results.

We can draw similar conclusions from Fig.6, which shows retrieval results on the CIRCO test dataset. In the first row, the proxy image provides a realistic and accurate depiction of the target image, allowing for improved retrieval results. In the second row, relying solely on the query image may cause the retrieval to overly focus on the beach background, contradicting the “no people” text information. Here, using the robust proxy (PI+RP) further enhances the overlooked aspects. In the third row, the balance metric balances the ‘the woman with sunglasses’ in proxy features with ‘trees’ in baseline features, achieving better retrieval results for both the subject and the background.

Ablation on Robust Features. We analyze the composition of our robust proxy on the CIRCO dataset, with results shown in Fig.5. The robust proxy is composed of three main components: q_i represents the query image features, s_i represents the semantic perturbation, and p_i represents our proxy image. A comparison among ①, ②, ③, and ④ shows that directly using our proxy image features can enhance retrieval performance to a certain extent, whereas relying solely on the query image or semantic perturbation tends to focus excessively on incorrect information, leading to a decline in performance. Additionally, comparisons among ⑤, ⑥, and ⑦ demonstrate that combining all three features together improves retrieval accuracy.

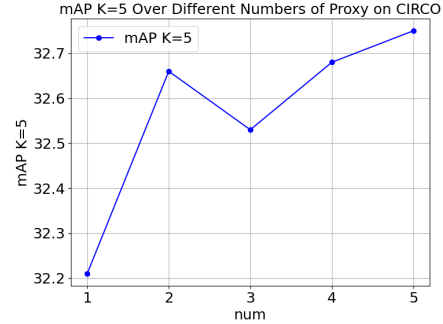


Figure 8. **Parameter Analysis.** The effect of varying numbers of proxy images on performance for the CIRCO dataset.

Effects of the weight λ . We analyze the λ used in Eq. (2). The experiments revealed that differences in dataset properties, result in different optimal values for λ . As shown in Fig.7, we found that the suitable λ for CIRCO, CIRR, and FashionIQ are around 0.3, 0.0, and 0.8.

Effects of the number of proxy images. We tested the impact of using different numbers of proxy images on performance with the CIRCO dataset. As shown in Fig.8, retrieval performance improves as the number of proxy images increases, though the rate of improvement gradually slows down. Therefore, the number of proxy images can be chosen based on a balance between efficiency and accuracy.

5. Conclusion

In this paper, we propose the IP-CIR, a plug-and-play training-free method for current ZSCIR tasks. Our method imagines proxy images, which align with the query image and the relative captions, providing fine-grained details that may be overlooked by the text. We combine the proxy images, query images, and semantic perturbation into more robust proxy image features, and propose a balance metric to compose the proxy and baseline similarities, enhancing the retrieval performance. The experiment results in CIRR, CIRCO, and FashionIQ datasets show that our method can

successfully improve retrieval accuracy and demonstrate the potential for providing additional retrieval information from the image side, revealing a new direction for improving retrieval accuracy in the future.

References

- [1] Lorenzo Agnolucci, Alberto Baldrati, Marco Bertini, and Alberto Del Bimbo. isearle: Improving textual inversion for zero-shot composed image retrieval. *arXiv preprint arXiv:2405.02951*, 2024. 1
- [2] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. 2, 3, 4, 6
- [3] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Qinsheng Zhang, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, Tero Karras, and Ming-Yu Liu. ediff-i: Text-to-image diffusion models with ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. 1
- [4] Alberto Baldrati, Marco Bertini, Tiberio Uricchio, and Alberto Del Bimbo. Effective conditioned and composed image retrieval combining clip-based features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21466–21474, 2022. 1
- [5] Alberto Baldrati, Marco Bertini, Tiberio Uricchio, and Alberto Del Bimbo. Conditioned and composed image retrieval combining and partially fine-tuning clip-based features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4959–4968, 2022. 1
- [6] Alberto Baldrati, Lorenzo Agnolucci, Marco Bertini, and Alberto Del Bimbo. Zero-shot composed image retrieval with textual inversion, 2023. 1, 2, 5
- [7] Alberto Baldrati, Lorenzo Agnolucci, Marco Bertini, and Alberto Del Bimbo. Zero-shot composed image retrieval with textual inversion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15338–15347, 2023. 1, 2, 6
- [8] Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models, 2023. 2
- [9] Xi Chen, Lianghua Huang, Yu Liu, Yujun Shen, Deli Zhao, and Hengshuang Zhao. Anydoor: Zero-shot object-level image customization. *arXiv preprint arXiv:2307.09481*, 2023. 2
- [10] Yanbei Chen and Loris Bazzani. Learning joint visual semantic matching embeddings for language-guided retrieval. page 136–152, Berlin, Heidelberg, 2020. Springer-Verlag. 1
- [11] Yanbei Chen, Shaogang Gong, and Loris Bazzani. Image search with text feedback by visiolinguistic attention learning. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2998–3008, 2020.
- [12] Yiyang Chen, Zhedong Zheng, Wei Ji, Leigang Qu, and Tat-Seng Chua. Composed image retrieval with text feedback via multi-grained uncertainty regularization, 2024.
- [13] Sanghyuk Chun. Improved probabilistic image-text representations, 2024.
- [14] Ginger Delmas, Rafael Sampaio de Rezende, Gabriela Csurka, and Diane Larlus. Artemis: Attention-based retrieval with text-explicit matching and implicit similarity, 2022.
- [15] Eric Dodds, Jack Culpepper, Simao Herdade, Yang Zhang, and Kofi Boakye. Modality-agnostic attention fusion for visual search with text feedback, 2020. 1
- [16] Aditya Ramesh et al. Hierarchical text-conditional image generation with clip latents, 2022. 1
- [17] Geonmo Gu, Sanghyuk Chun, Wonjae Kim, , Yoohoon Kang, and Sangdoo Yun. Language-only training of zero-shot composed image retrieval. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 1, 2, 6
- [18] Ryota Hinami and Shin’ichi Satoh. Discriminative learning of open-vocabulary object retrieval and localization by negative phrase augmentation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2605–2615, Brussels, Belgium, 2018. Association for Computational Linguistics. 1
- [19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1, 2
- [20] Shyamgopal Karthik, Karsten Roth, Massimiliano Mancini, and Zeynep Akata. Vision-by-language for training-free compositional image retrieval. *International Conference on Learning Representations (ICLR)*, 2024. 2, 6
- [21] Jooyeon Kim, Eulrang Cho, Sehyung Kim, and Hyunwoo J. Kim. Retrieval-augmented open-vocabulary object detection, 2024. 1
- [22] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, 2023. 3, 4
- [23] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. *CVPR*, 2023. 2
- [24] You Li, Fan Ma, and Yi Yang. Anysynth: Harnessing the power of image synthetic data generation for generalized vision-language tasks. *arXiv preprint arXiv:2411.16749*, 2024. 1
- [25] Chao Liang, Fan Ma, Linchao Zhu, Yingying Deng, and Yi Yang. Caphuman: Capture your moments in parallel universes. In *CVPR*, pages 6400–6409, 2024. 1
- [26] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context, 2015. 6

- [27] Zheyuan Liu, Cristian Rodriguez-Opazo, Damien Teney, and Stephen Gould. Image retrieval on real-life images with pre-trained vision-and-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2125–2134, 2021. [1](#), [2](#), [5](#)
- [28] Antoine Miech, Jean-Baptiste Alayrac, Ivan Laptev, Josef Sivic, and Andrew Zisserman. Thinking fast and slow: Efficient text-to-visual retrieval with transformers, 2021. [1](#)
- [29] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models, 2022. [1](#)
- [30] OpenAI. Gpt-4 technical report, 2023. [2](#), [3](#), [4](#)
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, A. Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021. [2](#)
- [32] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021. [1](#), [2](#)
- [33] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, Seyedeh Sara Mahdavi, Raphael Gontijo Lopes, Tim Salimans, Jonathan Ho, David Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. 2022. [2](#)
- [34] Kuniaki Saito, Kihyuk Sohn, Xiang Zhang, Chun-Liang Li, Chen-Yu Lee, Kate Saenko, and Tomas Pfister. Pic2word: Mapping pictures to words for zero-shot composed image retrieval. *CVPR*, 2023. [1](#), [2](#)
- [35] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv:2010.02502*, 2020. [1](#), [2](#)
- [36] Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. A corpus for reasoning about natural language grounded in photographs, 2019. [5](#)
- [37] Xudong Wang, Trevor Darrell, Sai Saketh Rambhatla, Rohit Girdhar, and Ishan Misra. Instancediffusion: Instance-level control for image generation, 2024. [2](#)
- [38] Zheng Wang, Xin Yuan, Toshihiko Yamasaki, Yutian Lin, Xin Xu, and Wenjun Zeng. Re-identification = retrieval + verification: Back to essence and forward with a new metric. *arXiv preprint arXiv:2011.11506*, 2020. [1](#)
- [39] Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. Elite: Encoding visual concepts into textual embeddings for customized text-to-image generation. *arXiv preprint arXiv:2302.13848*, 2023. [2](#), [4](#)
- [40] Hui Wu, Yupeng Gao, Xiaoxiao Guo, Ziad Al-Halah, Steven Rennie, Kristen Grauman, and Rogerio Feris. The fashion iq dataset: Retrieving images by combining side information and relative natural language feedback. *CVPR*, 2021. [1](#), [2](#), [5](#), [6](#)
- [41] Jinheng Xie, Yuexiang Li, Yawen Huang, Haozhe Liu, Wentian Zhang, Yefeng Zheng, and Mike Zheng Shou. Boxdiff: Text-to-image synthesis with training-free box-constrained diffusion. *arXiv preprint arXiv:2307.10816*, 2023. [2](#)
- [42] Zhengyuan Yang, Jianfeng Wang, Zhe Gan, Linjie Li, Kevin Lin, Chenfei Wu, Nan Duan, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. Reco: Region-controlled text-to-image generation. In *CVPR*, 2023. [2](#)
- [43] Zhenyu Yang, Dizhan Xue, Shengsheng Qian, Weiming Dong, and Changsheng Xu. Ldre: Llm-based divergent reasoning and ensemble for zero-shot composed image retrieval. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 80–90, 2024. [1](#), [2](#), [3](#), [4](#), [6](#)
- [44] Kai Zhang, Yi Luan, Hexiang Hu, Kenton Lee, Siyuan Qiao, Wenhui Chen, Yu Su, and Ming-Wei Chang. Magi-cLens: Self-supervised image retrieval with open-ended instructions. In *Proceedings of the 41st International Conference on Machine Learning*, pages 59403–59420. PMLR, 2024. [1](#), [2](#)
- [45] Liang Zheng, Yi Yang, and Alexander G. Hauptmann. Person re-identification: Past, present and future, 2016. [1](#)
- [46] Dewei Zhou, Zongxin Yang, and Yi Yang. Pyramid diffusion models for low-light image enhancement. *arXiv preprint arXiv:2305.10028*, 2023. [1](#)
- [47] Dewei Zhou, You Li, Fan Ma, Zongxin Yang, and Yi Yang. Migc: Multi-instance generation controller for text-to-image synthesis. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6818–6828, 2024. [2](#), [6](#)
- [48] Dewei Zhou, You Li, Fan Ma, Zongxin Yang, and Yi Yang. Migc++: Advanced multi-instance generation controller for image synthesis. *ArXiv*, abs/2407.02329, 2024. [2](#), [4](#)
- [49] Dewei Zhou, Ji Xie, Zongxin Yang, and Yi Yang. 3dis: Depth-driven decoupled instance synthesis for text-to-image generation. *ArXiv*, abs/2410.12669, 2024. [1](#)

Appendix

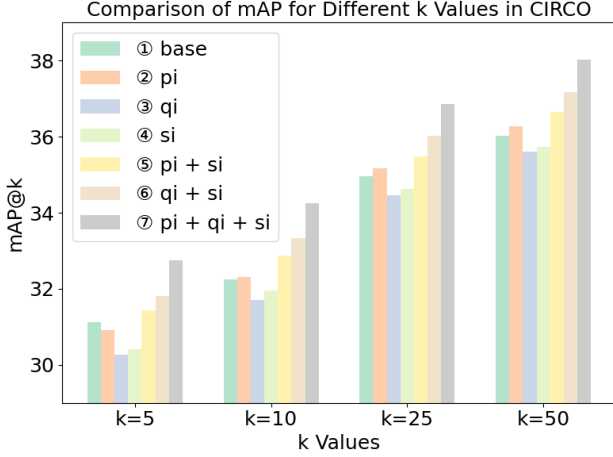


Figure 1. **Ablation results on the composition of robust proxy features** in the CIRCO dataset. **pi** indicates proxy features, **si** represents semantic perturbation, and **qi** indicates the query features.

A. More Analysis

In this section, we will conduct a deeper analysis of the experiments to better illustrate the details of our method.

A.1. The Construction of Robust Proxy

We analyze **the specific roles of each component in the robust features through Fig.1 and more visual examples in Fig.2**. The RP features include proxy image features, query image features, and semantic perturbation. Using the CIRCO dataset as an example, we construct RP features with different components for retrieval and obtain the top-1 retrieved image. This allows us to examine how varying the RP features influences the retrieved images.

We can draw the following conclusions: **1) The proxy image can provide richer information.** Since the proxy image contains semantic editing information as well as some information from the query image, using the proxy image can offer more comprehensive information. In Fig.1, compared to the baseline retrieval results ①, ② improves retrieval precision at $k = 10, 25, 50$, although it slightly reduces precision at $k = 5$. In contrast, directly using semantic perturbation or query features leads to an overall decline in retrieval precision. It can also be observed in Fig.2 that using only the query image (Qi) results in **an incorrect emphasis on background, even directly retrieving the original image itself**. While the second row shows that using only semantic perturbation (Si) **ignores details in the image such as angle**. In contrast, the proxy image achieves relatively better retrieval results because it provides background and angle information. **2) All the features are im-**

portant. In Fig.1, compared to ② and ③, when semantic perturbation is applied to enhance the text-driven editing information in proxy or query image features (as shown in ⑤ and ⑥), retrieval precision improves, highlighting the importance of semantic perturbation for robust features. Besides, the comparisons between ⑤ and ⑦, as well as ⑥ and ⑦, demonstrate that adding either proxy image features or query image features can further enhance the effectiveness of robust features. The second column in Fig.5 indicates that combining query image features can better preserve certain textual patterns and textures that are challenging for the proxy image to generate accurately.

A.2. The analysis of hyperparameters.

Since object content, associated text formats, and the alignment between text and target images vary across datasets, different datasets require different weighting parameters λ during retrieval. Additionally, the quality of proxies constructed for different datasets may also vary. For example, we observed that in FashionIQ, **the constructed proxies struggle to fully describe corresponding text for logo patterns and face challenges in generating pure white backgrounds during the generation process**. As a result, the role of proxy images in FashionIQ is relatively weaker, necessitating a larger λ . Moreover, if different backbones are used, the extracted features may emphasize different attributes, requiring dataset-specific adjustments to the λ parameter during retrieval.

Users can refine robust features by adjusting feature weights to emphasize specific retrieval aspects. For example, increasing the weight of the query image enhances details like logos or highlights attributes such as angle or style. Conversely, emphasizing proxy images or semantic perturbations shifts focus toward text-based editing directions.

B. More Qualitative Results

In this section, we show more results on three datasets. We present the Query image, relative caption, one of our generated Proxy images, as well as the top-1 baseline retrieval result and top-1 retrieval result of our method.

B.1. More Qualitative results on CIRCO.

We show qualitative results on CIRCO in Fig.3. We can draw the following conclusions: 1) The generated proxy images demonstrate certain detailed features, such as the yellow car in the first column, the desert and blue skateboard in the second column, the forest background and three horses in the third column, and the screen and keyboard in the last column. 2) Our method improves the Top-1 retrieval results to better meet the requirements to some extent. For exam-



Figure 2. **Ablation result on the Robust Proxy.** We present the visualization result of using different compositions (**Qi** represents only using query image, **Si** represents only semantic perturbation, and **Pi** represents only using the proxy image) of features in the robust proxy.

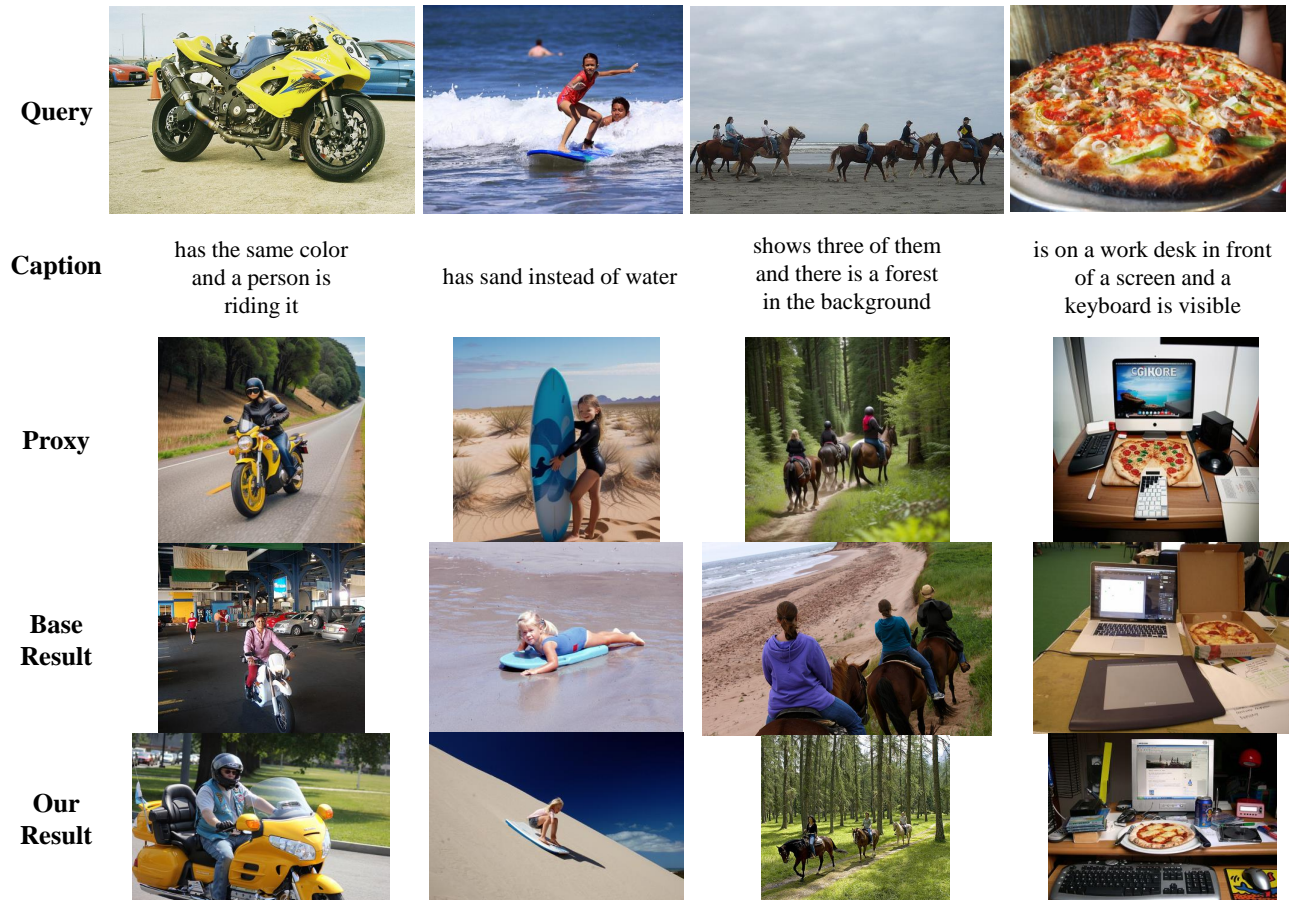


Figure 3. **Qualitative results on CIRCO dataset.** We show more improvement in top-1 retrieval results in the CIRCO dataset.

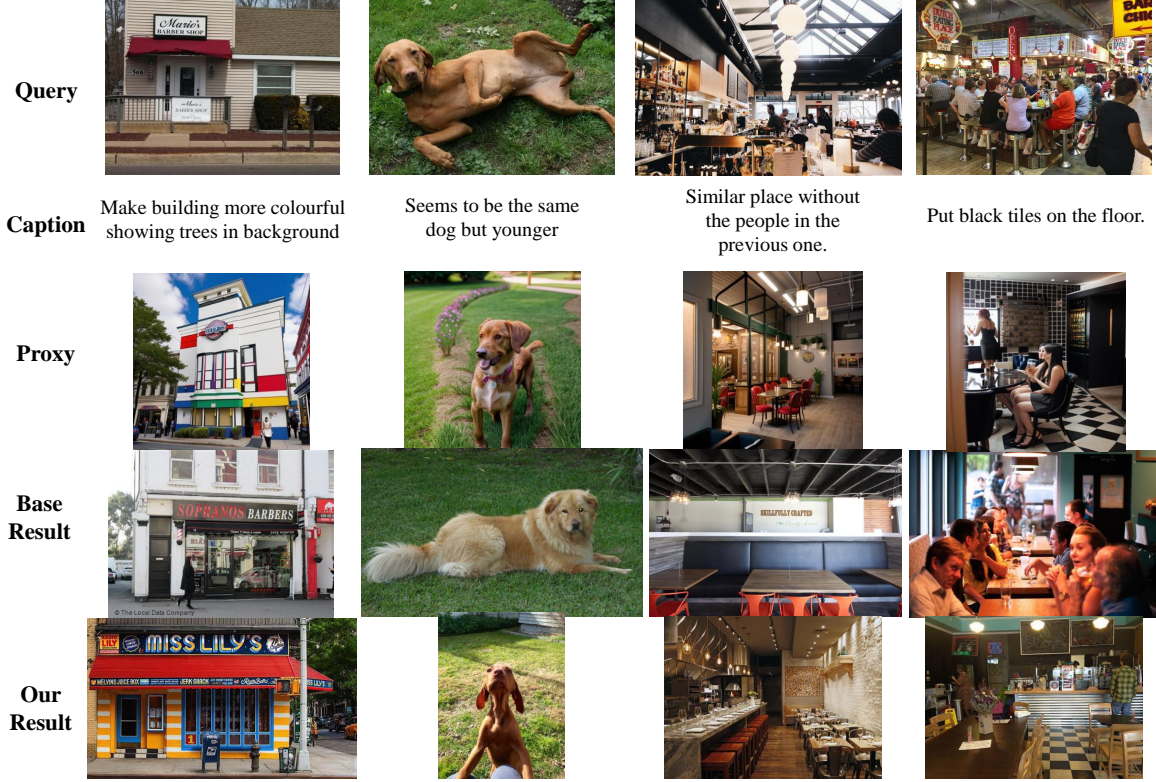


Figure 4. **Qualitative results on CIRR dataset.** We show more improvement in top-1 retrieval results in the CIRR dataset.

ple, in the first column, the car we retrieved is yellow. In the second column, the background of our result has sand instead of water. In the third column, with a proxy to imagine the forest scene, we successfully retrieved three horses in a forest. In the last column, the retrieved result closely matches the features presented by the proxy and better fits the described spatial characteristics.

B.2. More Qualitative results on CIRR.

We show more qualitative results on CIRR in Fig.4 with the improvements in TOP-1 retrieval performance. As shown, our generated proxy images provide elements such as the target’s background ambiance and scene (e.g., more colorful background and black tiles), as well as the type and fine-grained attributes of the main objects (e.g., the same type of dog). These enhancements contribute to improved retrieval accuracy.

B.3. More Qualitative results on FashionIQ.

We show more qualitative results on FashionIQ in Fig.5 with the improvements in TOP-1 retrieval performance. FashionIQ provides text descriptions of attributes such as color and patterns, enabling our proxy features to achieve relatively better retrieval results. At the same time, we believe that query image features are also important for the FashionIQ dataset. For example, in the second column, al-

though the baseline can retrieve clothes with pink color and black text, incorporating original image features allows the retrieval of images with logos that are more similar to the query while also aligning with the text description.

C. Limitation

Additional time overhead. While our method is plug-and-play in most scenarios and improves retrieval accuracy, it does introduce some time overhead. The process of layout generation and image generation adds certain time costs.

Sensitive to hyperparameters. The image-based retrieval enhancement is influenced by the performance of the constructed proxy images and the trade-off parameters used. Users need to carefully set reasonable weighting hyperparameters for the retrieval process. Thus, how to design a more reasonable method for combining weights or metrics to reduce sensitivity to parameters is an important direction for future research.





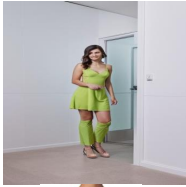

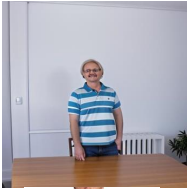









Query				
Caption	is a lime green	The shirt is pink in color with black writing.	is lighter blue with white stripes	its a greenish long dress
Proxy				
Base Result				
Our Result				

Figure 5. **Qualitative results on FashionIQ dataset.** We show more improvement in top-1 retrieval results in the FashionIQ dataset.