

# SORCE: Small Object Retrieval in Complex Environments

Chunxu Liu<sup>1,2\*</sup> Chi Xie<sup>2,3\*</sup> Xiaxu Chen<sup>2,4</sup> Wei Li<sup>2</sup> Feng Zhu<sup>2</sup> Rui Zhao<sup>2</sup> Limin Wang<sup>1,5†</sup>

<sup>1</sup>State Key Laboratory for Novel Software Technology, Nanjing University

<sup>2</sup>Sensetime Research <sup>3</sup>Tongji University <sup>4</sup>Beijing Institute of Technology <sup>5</sup>Shanghai AI Lab

<https://github.com/MCG-NJU/SORCE>

## Abstract

Text-to-Image Retrieval (T2IR) is a highly valuable task that aims to match a given textual query to images in a gallery. Existing benchmarks primarily focus on textual queries describing overall image semantics or foreground salient objects, possibly overlooking inconspicuous small objects, especially in complex environments. Such small object retrieval is crucial, as in real-world applications, the targets of interest are not always prominent in the image. Thus, we introduce **SORCE** (Small Object Retrieval in Complex Environments), a new subfield of T2IR, focusing on retrieving small objects in complex images with textual queries. We propose a new benchmark, SORCE-1K, consisting of images with complex environments and textual queries describing less conspicuous small objects with minimal contextual cues from other salient objects. Preliminary analysis on SORCE-1K finds that existing T2IR methods struggle to capture small objects and encode all the semantics into a single embedding, leading to poor retrieval performance on SORCE-1K.

Therefore, we propose to represent each image with multiple distinctive embeddings. We leverage Multimodal Large Language Models (MLLMs) to extract multiple embeddings for each image instructed by a set of Regional Prompts (ReP). Experimental results show that our multi-embedding approach through MLLM and ReP significantly outperforms existing T2IR methods on SORCE-1K. Our experiments validate the effectiveness of SORCE-1K for benchmarking SORCE performances, highlighting the potential of multi-embedding representation and text-customized MLLM features for addressing this task.

## 1 Introduction

While traditional Text-to-Image Retrieval (T2IR) focuses on retrieving images that match a given textual query, an important yet often overlooked challenge is retrieving small objects within complex scenes. In many real-world applications, the target of interest is not the dominant subject, but a small, inconspicuous object in a cluttered background. Despite significant advancements in T2IR methods [25, 38, 43, 39], existing frameworks struggle with such fine-grained retrieval, highlighting the need for dedicated research.

The challenge of retrieving small objects in complex environments with textual queries is highly relevant to various real-world applications. In video surveillance, it is essential for discovering inconspicuous yet critical items, such as concealed weapons in crowded areas or abandoned luggage in transportation hubs; in smart city systems, it helps detect traffic violations (e.g., unlicensed vehicles) and monitor infrastructure integrity (e.g., cracks in bridges). In cybersecurity, it can help discover

\*Equal Contribution. Work is done during internship at Sensetime.

†Corresponding author (lmwang@nju.edu.cn).



Figure 1: **Instances from different retrieval benchmarks.** We draw the bounding boxes of the referred objects in the image.

tampered watermarks and manipulated content; in consumer applications, it enables users to locate specific small objects within personal photo collections efficiently.

Despite the importance of small object retrieval, existing T2IR benchmarks primarily focus on retrieving images based on holistic descriptions rather than identifying fine-grained, non-salient objects. Popular T2IR benchmarks like Flickr30K [36] and COCO [15] highlights the main object or the whole picture in the image. Other datasets with image-text pairs [23, 2] provide more detailed descriptions, but still focus on describing the full image rather than specific small objects.

To address this gap, we introduce **Small Object Retrieval in Complex Environments (SORCE)**, a task focusing on retrieving the target image containing a specific small object. Here, *small* objects refer to the objects that are *visually* and *resolution-wise* small in the image, not semantically small.

For this task, we construct the SORCE-1K benchmark by collecting images from the SA-1B dataset [44]. As illustrated in Fig. 1, it is a high-resolution segmentation dataset with 1B mask annotations across 11M images, guaranteeing a high level of environment complexity. For each image, we identify a small target object and carefully craft a descriptive caption that uniquely specifies it, minimizing ambiguity within the benchmark. The resulting dataset, **SORCE-1K**, is a comprehensive benchmark for evaluating small object retrieval in complex environments.

In image retrieval, candidate images are typically pre-encoded into features, and retrieval is performed based on feature similarity with the query. Consequently, the image feature extraction is independent of the query content. This presents a significant challenge for feature extraction to capture every fine-grained detail within a single feature representation, making small object retrieval particularly difficult. Especially, the objects of interest in our SORCE-1K are often non-salient and embedded in complex scenes. Thus, relying on a single image feature for SORCE is likely unsatisfying.

To tackle the challenges of the SORCE task, we propose a simple yet effective starting point: representing an image with multiple features instead of a single embedding. Conventional feature extractors tend to prioritize the most salient objects, often overlooking fine-grained details. Extracting multiple feature representations can help mitigate this issue by capturing different aspects of an image. This can be naturally achieved using feature extractors based on Multimodal Large Language Models (MLLMs). Unlike traditional encoders, MLLMs can generate different feature embeddings based on different prompts [31, 8, 42], allowing them to focus on different aspects of the image while preserving global context. By leveraging this property, we can obtain a more comprehensive representation that improves retrieval performance for small, non-salient objects in complex environments.

Our main contributions are summarized as follows:

1. We introduce SORCE, a new subfield of Text-to-Image Retrieval (T2IR) that addresses the challenge of retrieving small, less conspicuous objects within images with complex environments based on textual queries.
2. We develop a new benchmark, SORCE-1K, for facilitating research on the SORCE task. SORCE-1K comprises 1,023 carefully curated images with complex backgrounds and textual queries that describe less prominent small objects with minimum surrounding context.
3. We propose an MLLM-based multi-embeddings approach, with Regional Prompts (ReP) to extract multiple image features focusing on different aspects. Superior performances

over existing T2IR methods on SORCE-1K highlight the potential of multi-embedding representation and text-guided MLLM features for addressing the SORCE task.

## 2 Related Work

### 2.1 Text-to-Image Retrieval Benchmarks

Text-to-Image Retrieval (T2IR) aims to retrieve relevant images from a candidate pool based on a given textual query. In practice, the candidate pool consists of pre-extracted image features, agnostic of the query content. The common approach is to first extract the text query embedding, then compute the cosine similarity between this text embedding and the embeddings of the candidate images. The top-ranked images with the highest similarity scores are retrieved as the final results.

Common benchmarks such as Flickr30K [36] and COCO [15] provide five brief captions for each image, describing the main object or event. Other datasets such as ShareGPT4V [2] and DOCCI [23] further extend the captions of images to be more detailed descriptions. Urban-1K [40] is another retrieval dataset for urban images with detailed descriptions. However, all of the above retrieval benchmarks are dedicated to enriching the query to be a detailed description of a whole image or the obvious objects in the image. In contrast, our benchmark SORCE-1K introduces a novel setting, which only contains the minimum required context in the query description. This approach aims to prevent retrieval models from taking shortcuts by recognizing and relying on non-target contextual descriptions. Therefore, our benchmark supports the robust retrieval evaluation of the specified target.

### 2.2 Visual Discovery of Small Objects

The visual discovery of small objects remains a longstanding challenge in computer vision, intersecting with several related domains, including the following. **I2I retrieval.** Handling small objects in I2I retrieval presents significant challenges, as highlighted in [28]. Prior works have addressed this by improving traditional or CNN-based descriptors [7, 3]. However, these methods are limited to single-modality retrieval and are not aligned with textual semantics. **Small object detection/segmentation.** Most small object detection or segmentation methods focus on a limited set of classes [13, 22, 35, 30, 27]. Open-vocabulary approaches [21, 41] still rely on a predefined set of target objects. Even advanced models like SAM [44] and DINO-X [26], which can detect objects without text prompts, often over-segment the image and are not suitable for our SORCE task. **Fine-grained VQA.** MLLMs [16, 17, 11] have been increasingly used to tackle fine-grained visual understanding. Benchmarks like SEED [10], MME [6], and MMBench [19] test models on tasks involving spatial reasoning, attribute comparison, object recognition, and spotting subtle differences. V\* [33] introduces a visual search task that requires grounding small objects in high-resolution images based on questions, pushing MLLMs to identify fine-grained details progressively. Unlike these VQA-based approaches, SORCE evaluates fine-grained perception through retrieval, where the model is not given any prior knowledge of the target object. This provides a complementary view of the model’s ability to localize and recognize small details without explicit textual input.

### 2.3 MLLM Embeddings

Recently, a new line of research has explored the use of MLLMs as feature extractors [9]. Trained with next-word prediction loss, MLLMs can effectively process and integrate visual and textual information. E5-V [9] first discovered that, when provided with carefully crafted prompts under the “one-word limitation” (e.g., “<image> Summarize the above image in one word:”), MLLMs leverage their multimodal comprehension to predict a one-word summary in the final token, effectively condensing input tokens into a compact embedding, and eliminating the feature gap between different modalities. Building on this insight, subsequent studies have shown that MLLMs have significant potential to overcome the limitations of conventional vision-language models (VLMs) and achieve outstanding performance in multimodal retrieval tasks [18, 42, 24, 20, 32].

In this work, we extend these findings by demonstrating that MLLMs can extract text-guided features through simple prompting. We highlight their potential for small object retrieval in complex scenes, leveraging their instruction-following ability.

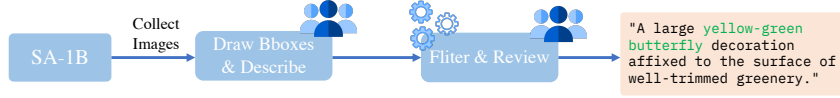


Figure 2: The construction process for SORCE-1K.

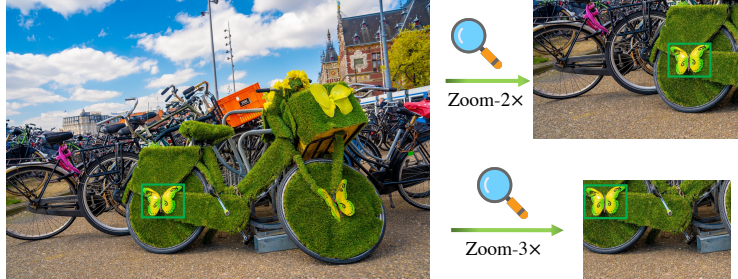


Figure 3: **Example of the difficulty levels.** We perform zooming  $2\times$  and  $3\times$  and crop the image while ensuring the object is in the resulting frame. More qualitative examples in Appendix D.

### 3 SORCE-1K Benchmark

#### 3.1 Dataset Highlights and Statistics

The proposed benchmark follows the same task formulation and evaluation protocol as existing text-to-image retrieval benchmarks. However, its key challenges differ from traditional retrieval benchmarks, as it emphasizes retrieving the image that contains a specific small target in potentially complex environments. Below, we highlight the key characteristics of this benchmark.

**Small targets.** SORCE-1K requires a model to retrieve an image containing a specific small target object. The target is typically a less salient element within the image, often blending into the background or hidden among details. This presents a greater challenge compared to traditional retrieval settings, where the target is usually prominent foreground objects or the whole scene.

**Complex environments.** The targets are usually embedded in complex scenes, often appearing as background components or fine details. The scene complexity is ensured by using SA-1B [44] as the image source and selecting only images with a large number of masks. This requires the model to recognize and differentiate between entire scenes and the objects within them. Note that the numbers of masks in images in our benchmark are presented in Appendix A.

**Local and detailed descriptions.** Each target is referenced by a detailed description that uniquely matches it without overlapping with objects in other images from the dataset. The corresponding description primarily focuses on local features of the target object, rather than its surrounding context or the entire scene. This ensures that the retrieval model need to identify the target itself rather than relying on contextual cues mentioned in the descriptions as shortcuts.

**Multiple difficulty levels.** Retrieving full-size images can be too challenging. To facilitate comparison and analysis across different scales, SORCE-1K is designed with two additional “zoom-in” settings where the target object is kept while the image context is cropped. The 3 settings (Full Res., Zoom- $2\times$ , Zoom- $3\times$ ) are depicted in Fig. 3. This helps understand gaps between traditional retrieval tasks and the proposed task setting.

**High quality.** The entire benchmark is annotated by researchers themselves rather than crowd-sourced annotators. Additionally, measures such as manual checking are implemented during the annotation process to ensure uncompromised quality.

**Dataset statistics.** To construct the dataset, a total of 2.1K images from SA-1B are labeled and examined, and about half are discarded for not meeting these criteria. The result dataset is an evaluation-only benchmark with 1,023 images. Each image is annotated with one description and one box, though the box is not used for the retrieval task. The description refers to the object in the box on the image. As is shown in Fig. 4(a), the descriptions in the dataset are at least 6 words long, at most 42 words, with an average of 16.9 words. The bounding boxes are very small compared with

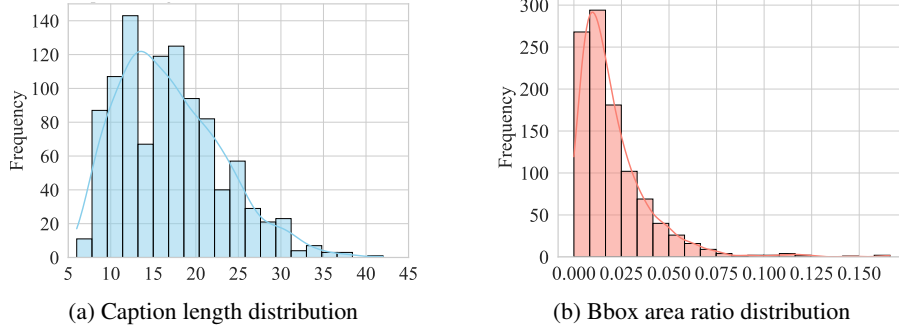


Figure 4: **Statistics for the proposed SORCE-1K benchmark.**

the high-resolution image. As is shown in Fig. 4 (b), the small objects to retrieve usually take up less than 10% of the images. The evaluation metric is the same as standard retrieval benchmarks.

### 3.2 Dataset Construction

This benchmark is annotated manually by researchers, followed by automatic and manual verification.

**Data source.** To construct a benchmark featuring small objects in complex environments, we use images from SA-1B [44]. It has (1) high-resolution images, which are more likely to contain small objects, and (2) complex scenes, quantifiable by the number of masks in SA-1B annotations.

**Region-description pair annotation.** For each image, we manually check whether it contains a target that (1) is relatively small compared to the entire image, (2) can be described using natural language, and (3) has a unique combination of local features that distinguishes it from all other instances in the dataset. If such a target exists, we annotate it with a bounding box and a corresponding unique description; otherwise, we exclude the image from the dataset.

**Dataset tiering with difficulty.** For each Full Res. sample, we divide it by  $2 \times 2$  and  $3 \times 3$  to obtain a Zoom-2 $\times$  sample and a Zoom-3 $\times$  sample. To ensure the target object remains intact, we select the sub-image that covers the bounding box the most, and extend the region when necessary to fully include the bounding box. The description annotation is shared across all three levels. Compared to the Full Res. image, the Zoom-2 $\times$  setting is easier, as the target is more prominent, while the Zoom-3 $\times$  setting is the easiest.

**Automatic check.** We apply a few simple automatic filtering steps to the candidate samples from the previous step. These steps include (1) removing samples where the bounding box is too large and dominant in the image and (2) removing or correcting samples with overly short descriptions, as they are likely too simple or non-unique. These steps help ensure that annotation errors are minimized.

**Manual review.** Finally, we perform a manual review of the annotated dataset. Each image is visualized along with its annotated region-description pair. We examine the visualizations and remove samples where the descriptions are inaccurate or ambiguous, particularly those that may refer to multiple instances in the dataset. Please see Appendix C for more details regarding the steps above.

### 3.3 Other Application Scenarios

**Visual grounding and described object detection.** Since each image is annotated with a single bounding box associated with a language description, the proposed benchmark can serve directly as a visual grounding [37] benchmark. It is also applicable to described object detection [34], a related task that involves detecting all possible objects in an image set based on a given language description.

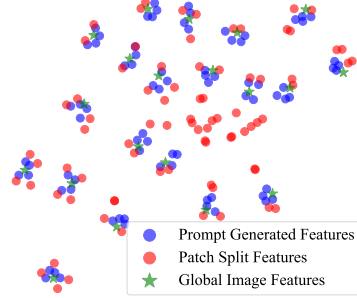
**Visual search.** Additionally, because the target object in each image is relatively small, the benchmark can be adapted for the visual search [33] task, where a model must answer a question about a small target within an image. This can be achieved by prompting an LLM to convert the target object’s description into a question-answer pair focused on a specific attribute of the object.

**Instance-level retrieval.** While the proposed benchmark is primarily for image-level retrieval, it is also applicable to instance-level retrieval tasks [1]. These tasks typically require explicit bounding box predictions, which our benchmark does not inherently demand but provides in the annotations.



(a) **Similarity matrix of regional prompts generated features ( $f_{A,B,C,D}$ ) and independently generated features from  $2 \times 2$  Split ( $F_{A,B,C,D}$ ).** Both sets of features are extracted using E5-V.  $F_{A,B,C,D}$  are extracted by prompt: “Summarize the above image in one word.”

	$f_A$	$f_B$	$f_C$	$f_D$
$F_{Global}$	0.85	0.84	0.85	0.88
$F_{Global}$	0.63	0.54	0.84	0.81
	$F_A$	$F_B$	$F_C$	$F_D$



(b) **Similarity matrix of global image feature ( $F_{global}$ ) and  $f_{A,B,C,D}$ ,  $F_{A,B,C,D}$ .**

(c) **Independently generated features from  $2 \times 2$  splits.** We randomly select 20 image samples from SORCE-1K.

Figure 5: **Qualitative comparisons with regional prompt generated features and independently generated features from  $2 \times 2$  splits.**

## 4 Method

### 4.1 Multiple Feature Representation

As shown in Tab. 1, existing approaches obtain poor results on the Full Res. setting of SCORCE-1K. This shows that compressing all visual detail into a single image feature remains highly challenging for existing approaches either based on VLMs [25, 29] or MLLMs [11, 9].

Here we propose to extract multiple features tailored to different aspects of an image, leveraging the instruction-following abilities of MLLMs. We observed that Regional Prompts (ReP), such as “<image> Summarize the [regional part] of the image in one word:”, tend to focus on summarizing specific regions of the image.

To show this, we begin by qualitatively analyzing Fig. 5a, which compares the cosine similarity between features generated using regional prompts ( $f_A, f_B, f_C, f_D$ ) and features independently generated from A, B, C, D splits ( $F_A, F_B, F_C, F_D$ ). The results reveal that  $f$  closely aligns with  $F$  for corresponding regions. Sequentially, in Fig. 5b, we compare  $f_{A,B,C,D}$  and  $F_{A,B,C,D}$  with the global image feature  $F_{Global}$ , showing that  $f_{A,B,C,D}$  is consistently closer to the global image feature than  $F_{A,B,C,D}$ . Taking this analysis further, we randomly select 20 images from our benchmark and visualize a t-SNE map of independently generated patch split features and regional prompt-generated features. Fig. 5c demonstrates that regional prompt-generated features cluster closely around the global image feature, while patch split features are often farther away, sometimes at significant distances. These observations demonstrate that regional prompt-generated features correspond to independently generated features of their respective regions, while retaining global information.

### 4.2 Contrastive Finetuning

To enhance the discriminability of MLLM embeddings, previous works [9, 18, 14, 42, 24] adopt contrastive fine-tuning widely. Following this, we provide a simple but effective solution, by integrating Regional Prompts into the fine-tuning process.

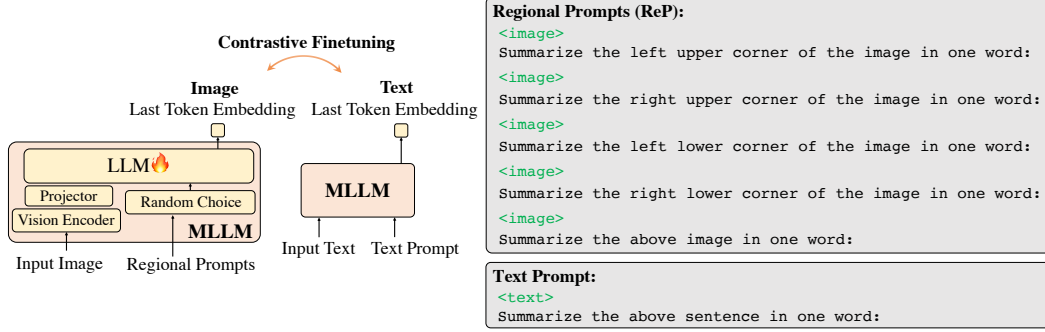


Figure 6: **Pipeline of our contrastive finetuning process.** *Random Choice* block means that we randomly choose a regional prompt from the right, and use the resulting image embedding to align with the text description for the corresponding regional regions.

To employ regional prompts in fine-tuning, we restructure the dataset to include descriptions of each region. Using InternVL2.5-38B [4], we recaption the dataset by prompting it (see Appendix E) to describe specific parts of the input image (*left upper*, *right upper*, *left lower* and *right lower*). This allows us to align prompt-generated features ( $f_{lu}, f_{ru}, f_{ll}, f_{rl}$ ) with their corresponding text features ( $t_{lu}, t_{ru}, t_{ll}, t_{rl}$ ). The text features are extracted using the prompt: “<text> Summarize the above sentence in one word:”. To prevent redundancy from the same image appearing multiple times, we randomly select one regional prompt per image in training process. The training objective is as follows:

$$L = \frac{1}{2} (L_{CE}(f, t) + L_{CE}(t, f)), \quad (1)$$

where  $L_{CE}$  is cross entropy loss and  $L_{CE}(a, b)$  is:

$$L_{CE}(a, b) = -\log \frac{\exp(\cos(a^i, b^i)/\tau)}{\sum_j \exp(\cos(a^i, b^j)/\tau)}. \quad (2)$$

$\cos(\cdot, \cdot)$  denotes the cosine similarity function and  $\tau$  is the temperature hyper-parameter.

### 4.3 Inference Setting

Since retrieving a small target in a complex scene is agnostic to the query content, the generated image features must encompass as much information as possible to facilitate effective querying. Leveraging the text-guided feature extraction capability of MLLMs, we propose to use multiple features rather than a single feature to represent each image. Specifically, for each image, we generate five distinct features using five regional prompts, as listed in Fig. 6. For a given query, we compute the cosine similarity between the query and each of the five features, selecting the closest one as the final feature for metric computation.

## 5 Experiments

### 5.1 Experiment Settings

**Datasets and benchmarks.** We use COCO-118K [15, 11] for MLLM fine-tuning, which consists of 118K image-text pairs. We recaption the dataset with InternVL2.5-38B [4], with detailed prompting instructions provided in Appendix E. After recaptioning, each image in COCO-118K is paired with four regional descriptions (*left upper corner*, *right upper corner*, *left lower corner*, *right lower corner*) and a summary caption.

For evaluation, we mainly evaluate our method on commonly used retrieval benchmarks, Flickr30K [36] and COCO [15] for validating the versatility of the proposed method in full image retrieval scenarios. In addition, we use the proposed benchmark to investigate the small object retrieval ability in complex environments, and set different difficulty levels, Zoom-3 $\times$ , Zoom-2 $\times$  and Full Res., with increasing levels of challenge.

Table 1: **Image retrieval results on SORCE-1K**. We report R@1, R@5, and R@10. The results indicate that ReP enhances the performance of MLLM-based feature extractors, while fine-tuning further improves retrieval recall. ft. means fine-tuned. Please refer to Sec. 4.2.

Method	Multiple Features	SORCE-1K (Ours)								
		Zoom-3×			Zoom-2×			Full Res.		
		R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
CLIP ViT-B [25]	✗	42.3	61.5	69.0	31.3	50.2	57.6	17.9	32.1	39.4
CLIP ViT-L [25]	✗	47.8	66.7	73.2	36.5	54.1	62.5	19.5	33.7	41.5
EVA-02-CLIP 5B [29]	✗	64.8	78.8	84.1	44.3	60.8	69.2	23.0	36.7	44.8
LLaVA-Next-8B [12]	✗	41.0	59.1	69.1	28.2	44.6	52.4	13.4	23.0	30.4
E5-V [9]	✗	57.6	74.0	82.2	42.4	62.2	69.0	21.9	36.7	44.6
E5-V + ReP	✓	62.0	80.6	86.3	54.0	73.2	80.2	27.7	45.4	53.7
E5-V (ft.) + ReP	✓	<b>68.0</b>	<b>85.5</b>	<b>90.9</b>	<b>56.3</b>	<b>77.1</b>	<b>83.0</b>	<b>31.5</b>	<b>50.6</b>	<b>60.0</b>

Table 2: **Image retrieval results on Flickr30K [36] and COCO [15]**. We also list text retrieval results for reference. The results indicate that MLLM-based feature extractors have comparable performance with their baseline model. ReP represents Regional Prompts.

Method	Multiple Features	Image Retrieval						Text Retrieval					
		Flickr30K			COCO			Flickr30K			COCO		
		R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
CLIP ViT-B [25]	✗	62.1	85.6	91.8	33.1	58.4	69.1	81.9	96.2	98.8	52.5	76.7	84.6
CLIP ViT-L [25]	✗	67.3	89.0	93.3	37.1	61.6	71.5	87.4	98.3	99.3	57.9	81.2	87.8
LLaVA-Next-8B [12]	✗	59.9	83.5	90.2	34.5	60.2	70.9	69.4	90.7	95.0	41.9	67.1	77.1
EVA-02-CLIP 5B [29]	✗	78.8	94.2	96.8	51.1	75.0	82.7	<b>93.9</b>	<b>99.4</b>	<b>99.8</b>	<b>68.8</b>	<b>87.8</b>	<b>92.8</b>
E5-V [9]	✗	80.8	95.5	97.7	<b>52.1</b>	<b>76.6</b>	83.6	88.1	98.8	99.4	62.2	83.6	89.9
LLaVA-Next-8B + ReP	✓	59.6	83.4	90.2	34.3	60.1	70.8	68.8	90.6	94.9	41.1	66.4	76.6
E5-V + ReP	✓	<b>81.0</b>	<b>95.7</b>	<b>97.8</b>	51.8	76.3	<b>84.4</b>	87.9	98.7	99.4	59.9	82.0	88.6

**Implementation details.** Our method is trained on E5-V, which is pretrained on LLaVA-Next-8B [12]. From another perspective, our method can be seen as a two-stage training, which is first pretrained on text-only datasets, and then finetuned on text-image pairs, sharing the same intuition with the training strategy of LamRA [18]. We fine-tune the LLM part of the MLLM by QLoRA [5], while keeping the image encoder and image feature projector frozen. We trained the model for a single epoch using 32 V100 GPUs, a batch size of 768, and a learning rate of  $2 \times 10^{-4}$ , requiring approximately 3.5 hours.

## 5.2 Comparison with Other Methods

The baselines to compare include CLIP with ViT-B and ViT-L-336px, and LLaVA-NeXT-8B, E5-V in MLLMs. For MLLMs, we use “<text> Summarize the above sentence in one word:” as the prompt for text feature extraction. For the proposed method, we extract five prompt-generated features for each image, while for other methods, each image is represented with one image feature.

We report Recall@K (R@K) with K=1, 5, 10 in Tab. 2 for image and text retrieval. For our benchmark in Tab. 1, we report image retrieval R@K with K=1, 5, 10. Compared to the baselines, our method keeps comparable performance on COCO and Flickr. And our method achieves competitive performance in our benchmark. Combining the performance on Flickr and COCO, our method showcases the potential of text-guided feature representation.

## 5.3 Ablation Study

**Does simple ensembling achieve the same performance?** In the proposed method, we leverage multiple features, each focusing on different aspects of the image, to better capture its fine-grained details. A natural question arises: Does the performance improvement stem solely from ensembling multiple MLLM features? To investigate this, we use four additional synonyms for “Summarize” as global feature extraction prompts: “Conclude”, “Synopsisize”, “Condense”, and “Encapsulate”. As shown in Tab. 3, these five prompts yield negligible improvement over the single-feature setting. This demonstrates that our multiple regional prompts genuinely alter the focus of the MLLM-generated image features, rather than merely benefiting from ensembling.

**Why not obtain multiple features by cropping?** Cropping is an intuitive approach to generating multiple features from an image. However, it has significant limitations: it lacks global context when extracting features from each cropped region and risks splitting the target object entirely. Here we

Table 3: Performance comparison with synonym prompts ensembling.

Methods	Multiple Features	SORCE-1K (R@5)		
		Zoom-3 $\times$	Zoom-2 $\times$	Full Res.
E5-V	$\times$	74.0	62.2	36.7
+ Synonym Prompts	$\checkmark$	74.5	61.8	36.5
+ ReP	$\checkmark$	<b>85.5</b>	<b>77.1</b>	<b>50.6</b>

Table 4: Performance comparison between  $2 \times 2$  split features and regional prompt features on SORCE-1K.

Methods	Multiple Features	SORCE-1K (R@5)		
		Zoom 3 $\times$	Zoom 2 $\times$	Full Res.
E5-V	$\times$	74.0	62.2	36.7
+ $2 \times 2$ Split	$\checkmark$	79.4	70.8	47.9
+ ReP	$\checkmark$	80.6	73.2	45.4
(ft.) + ReP	$\checkmark$	<b>85.5</b>	<b>77.1</b>	<b>50.6</b>

Table 6: Effects of Different Prompt Numbers. In sequence,  $P_a, P_b, P_c, P_d, P_e$  represent four regional prompts followed by one summary prompt.

Methods	SORCE-1K (R@5)		
	Zoom-3 $\times$	Zoom-2 $\times$	Full Res.
E5-V	74.0	62.2	36.7
+ $P_a$	64.4	54.0	32.0
+ $P_{a,b}$	65.9	56.7	36.5
+ $P_{a,b,c}$	73.6	67.7	39.3
+ $P_{a,b,c,d}$	75.2	69.8	42.2
+ $P_{a,b,c,d,e}$	<b>80.6</b>	<b>73.2</b>	<b>45.4</b>

Table 5: Performance comparison between  $2 \times 2$  split features and regional prompt features on Flickr and COCO.

Methods	Multiple Features	Flickr30K R@5	COCO R@5
E5-V	$\times$	95.5	<b>76.6</b>
+ $2 \times 2$ Split	$\checkmark$	84.8	59.4
+ ReP	$\checkmark$	<b>95.7</b>	76.3

Table 7: Comparison between directly contrastive finetuning and random choice finetuning. M-All means merge the regional descriptions of each recaptioned image into one caption, and perform contrastive fine-tuning.

Methods	Multiple Features	SORCE-1K (R@5)		
		Zoom 3- $\times$	Zoom 2- $\times$	Full Res.
E5-V	$\times$	74.0	62.2	36.7
+ ReP	$\checkmark$	80.6	73.2	45.4
(M-All ft.) + ReP	$\checkmark$	80.3	69.9	43.1
(ft.) + ReP	$\checkmark$	<b>85.5</b>	<b>77.1</b>	<b>50.6</b>

crop the image into  $2 \times 2$  non-overlapping splits and extract features from each region. We also append the global feature when evaluating, resulting in a set of five features per image and ensuring a fair comparison. In Tab. 4, cropping improves performance over the base model E5-V and surpasses regional prompts. However, this is largely due to our benchmark only involving minimal required context for identifying the target small object. In contrast, Tab. 5 shows that cropping significantly degrades performance, indicating that it is not a universally applicable solution.

**Does the number of the prompts affect performance?** We denote the left upper corner, right upper corner, left lower corner, right lower corner and global summary as  $a, b, c, d, e$  accordingly. Then we ablate with 1 to 5 prompts, from  $a$  to  $a, b, c, d, e$ , denoting as  $P_a, P_{a,b}, \dots, P_{a,b,c,d,e}$ . From Tab. 6, we can see that with more prompt-generated features included, the performance keeps improving.

**Does the regional prompt finetuning truly improve performance?** We transform our InternVL2.5-38B recaptioned COCO-118K dataset into a format where each image is associated with a single unified description, and apply the contrastive tuning between the global image feature and the full description. As shown in Tab. 7, the instruction-following capability exhibits a slight decline after fine-tuning, suggesting that our random choice fine-tuning strategy provides a beneficial effect.

We present more ablation study like the effect of using different prompts in Appendix F.

## 6 Conclusion

In this work, we focus on an overlooked demand in T2IR, which is retrieving a specific small object from images with complex scenes. We start by formulating the corresponding task setting, Small Object Retrieval in Complex Environments, and constructing the SORCE-1K benchmark. The benchmark focuses on retrieving small objects from large images, based on language descriptions, without the interference of contexts. Noticing that a single feature may not represent all small objects in an image, we leverage MLLMs to extract text-customized features focusing on various aspects in an image, which provides a simple and intuitive solution for this setting. We hope the SORCE-1K benchmark, together with the proposed solution, can serve as a foundation for future research on such a setting. We discuss the limitations and broader impact in Appendices G and H.

## References

- [1] Zhaowei Cai, Gukyeon Kwon, Avinash Ravichandran, Erhan Bas, Zhuowen Tu, Rahul Bhotika, and Stefano Soatto. X-detr: A versatile architecture for instance-wise vision-language tasks. In *European Conference on Computer Vision*, pages 290–308. Springer, 2022.
- [2] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. In *European Conference on Computer Vision*, pages 370–387. Springer, 2024.
- [3] Zhenfang Chen, Zhanghui Kuang, Wayne Zhang, and Kwan-Yee K Wong. Learning local similarity with spatial relations for object retrieval. In *Proceedings of the 27th ACM international conference on Multimedia*, pages 1703–1711, 2019.
- [4] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024.
- [5] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *Advances in neural information processing systems*, 36:10088–10115, 2023.
- [6] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.
- [7] Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, Teddy Furon, and Ondrej Chum. Efficient diffusion on region manifolds: Recovering small objects with compact cnn representations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2077–2086, 2017.
- [8] Ting Jiang, Shaohan Huang, Zhongzhi Luan, Deqing Wang, and Fuzhen Zhuang. Scaling sentence embeddings with large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 3182–3196, 2024.
- [9] Ting Jiang, Minghui Song, Zihan Zhang, Haizhen Huang, Weiwei Deng, Feng Sun, Qi Zhang, Deqing Wang, and Fuzhen Zhuang. E5-v: Universal embeddings with multimodal large language models. *arXiv preprint arXiv:2407.12580*, 2024.
- [10] Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. Seed-bench: Benchmarking multimodal large language models. In *CVPR*, pages 13299–13308, 2024.
- [11] Bo Li, Hao Zhang, Kaichen Zhang, Dong Guo, Yuanhan Zhang, Renrui Zhang, Feng Li, Ziwei Liu, and Chunyuan Li. Llava-next: What else influences visual instruction tuning beyond data?, 2024.
- [12] Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and Chunyuan Li. Llava-next: Stronger llms supercharge multimodal capabilities in the wild, 2024.
- [13] Jianan Li, Xiaodan Liang, Yunchao Wei, Tingfa Xu, Jiashi Feng, and Shuicheng Yan. Perceptual generative adversarial networks for small object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1222–1230, 2017.
- [14] Sheng-Chieh Lin, Chankyu Lee, Mohammad Shouybi, Jimmy Lin, Bryan Catanzaro, and Wei Ping. MM-EMBED: UNIVERSAL MULTIMODAL RETRIEVAL WITH MULTIMODAL LLMS. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [15] 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 *Computer vision—ECCV 2014: 13th European conference, zurich, Switzerland, September 6–12, 2014, proceedings, part v 13*, pages 740–755. Springer, 2014.

- [16] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023.
- [17] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024.
- [18] Yikun Liu, Pingan Chen, Jiayin Cai, Xiaolong Jiang, Yao Hu, Jiangchao Yao, Yanfeng Wang, and Weidi Xie. Lamra: Large multimodal model as your advanced retrieval assistant. *arXiv preprint arXiv:2412.01720*, 2024.
- [19] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? In *European conference on computer vision*, pages 216–233. Springer, 2024.
- [20] 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*, pages 2125–2134, 2021.
- [21] Matthias Minderer, Alexey Gritsenko, and Neil Houlsby. Scaling open-vocabulary object detection. *Advances in Neural Information Processing Systems*, 36:72983–73007, 2023.
- [22] Junhyug Noh, Wonho Bae, Wonhee Lee, Jinhwan Seo, and Gunhee Kim. Better to follow, follow to be better: Towards precise supervision of feature super-resolution for small object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9725–9734, 2019.
- [23] Yasumasa Onoe, Sunayana Rane, Zachary Berger, Yonatan Bitton, Jaemin Cho, Roopal Garg, Alexander Ku, Zarana Parekh, Jordi Pont-Tuset, Garrett Tanzer, et al. Docci: Descriptions of connected and contrasting images. In *European Conference on Computer Vision*, pages 291–309. Springer, 2024.
- [24] Yassine Ouali, Adrian Bulat, Alexandros Xenos, Anestis Zaganidis, Ioannis Maniadis Metaxas, Georgios Tzimiropoulos, and Brais Martinez. Discriminative fine-tuning of lvlms. *arXiv preprint arXiv:2412.04378*, 2024.
- [25] 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*, pages 8748–8763. PmLR, 2021.
- [26] Tianhe Ren, Yihao Chen, Qing Jiang, Zhaoyang Zeng, Yuda Xiong, Wenlong Liu, Zhengyu Ma, Junyi Shen, Yuan Gao, Xiaoke Jiang, Xingyu Chen, Zhuheng Song, Yuhong Zhang, Hongjie Huang, Han Gao, Shilong Liu, Hao Zhang, Feng Li, Kent Yu, and Lei Zhang. Dino-x: A unified vision model for open-world object detection and understanding, 2024.
- [27] Shengtian Sang, Yuyin Zhou, Md Tauhidul Islam, and Lei Xing. Small-object sensitive segmentation using across feature map attention. *IEEE transactions on pattern analysis and machine intelligence*, 45(5):6289–6306, 2022.
- [28] Torsten Sattler, Tobias Weyand, Bastian Leibe, and Leif Kobbelt. Image retrieval for image-based localization revisited. In *BMVC*, page 4, 2012.
- [29] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training techniques for clip at scale. *arXiv preprint arXiv:2303.15389*, 2023.
- [30] Huan Wang, Luping Zhou, and Lei Wang. Miss detection vs. false alarm: Adversarial learning for small object segmentation in infrared images. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8509–8518, 2019.
- [31] Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. Improving text embeddings with large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11897–11916, 2024.

- [32] Cong Wei, Yang Chen, Haonan Chen, Hexiang Hu, Ge Zhang, Jie Fu, Alan Ritter, and Wenhui Chen. Uniir: Training and benchmarking universal multimodal information retrievers. In *European Conference on Computer Vision*, pages 387–404. Springer, 2024.
- [33] Penghao Wu and Saining Xie. V\*: Guided visual search as a core mechanism in multimodal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13084–13094, 2024.
- [34] Chi Xie, Zhao Zhang, Yixuan Wu, Feng Zhu, Rui Zhao, and Shuang Liang. Described object detection: Liberating object detection with flexible expressions. *Advances in Neural Information Processing Systems*, 36:79095–79107, 2023.
- [35] Chenhongyi Yang, Zehao Huang, and Naiyan Wang. Querydet: Cascaded sparse query for accelerating high-resolution small object detection. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, pages 13668–13677, 2022.
- [36] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the association for computational linguistics*, 2:67–78, 2014.
- [37] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14*, pages 69–85. Springer, 2016.
- [38] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11975–11986, 2023.
- [39] Beichen Zhang, Pan Zhang, Xiaoyi Dong, Yuhang Zang, and Jiaqi Wang. Long-clip: Unlocking the long-text capability of clip. In *European Conference on Computer Vision*, pages 310–325. Springer, 2024.
- [40] Beichen Zhang, Pan Zhang, Xiaoyi Dong, Yuhang Zang, and Jiaqi Wang. Long-clip: Unlocking the long-text capability of clip. In *European Conference on Computer Vision*, pages 310–325. Springer, 2024.
- [41] Hao Zhang, Feng Li, Xueyan Zou, Shilong Liu, Chunyuan Li, Jianwei Yang, and Lei Zhang. A simple framework for open-vocabulary segmentation and detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1020–1031, 2023.
- [42] Xin Zhang, Yanzhao Zhang, Wen Xie, Mingxin Li, Ziqi Dai, Dingkun Long, Pengjun Xie, Meishan Zhang, Wenjie Li, and Min Zhang. Gme: Improving universal multimodal retrieval by multimodal llms. *arXiv preprint arXiv:2412.16855*, 2024.
- [43] Junjie Zhou, Zheng Liu, Shitao Xiao, Bo Zhao, and Yongping Xiong. Vista: Visualized text embedding for universal multi-modal retrieval. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3185–3200, 2024.
- [44] Xueyan Zou, Jianwei Yang, Hao Zhang, Feng Li, Linjie Li, Jianfeng Wang, Lijuan Wang, Jianfeng Gao, and Yong Jae Lee. Segment everything everywhere all at once. *Advances in neural information processing systems*, 36:19769–19782, 2023.

## Appendix

### A Environment Complexity of SORCE-1K

Since our SORCE-1K benchmark is collected from SA-1B [44], every image has the corresponding mask label. The complexity of the image usually increases with the number of segmentation masks. Therefore, we provide the segmentation mask number distribution of SORCE-1K in Fig. 7.

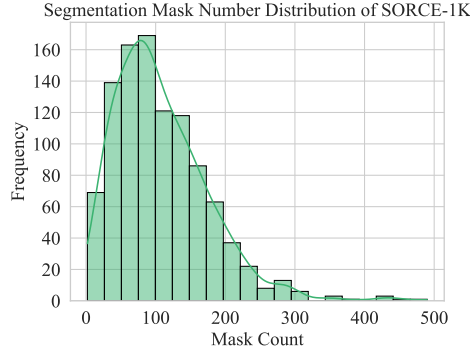


Figure 7: **Mask distribution of SORCE-1K.** Environment complexity usually grows with the number of segmentation masks.

### B The Caption Length Distribution of Flickr30K and COCO

We present the caption length distribution charts in Fig. 8, alongside the SORCE-1K caption length distribution in Fig. 4. As shown, the caption lengths in SORCE-1K are comparable to or slightly longer than those in COCO and Flickr. However, the key distinction lies in the focus of the captions: SORCE-1K emphasizes detailed descriptions of the target object while minimizing the inclusion of global context. Consequently, caption length is not the primary metric we aim to highlight.

### C Construction Details of SORCE-1K

Here we provide more details regarding the construction process (see Sec. 3) of the proposed SORCE-1K benchmark.

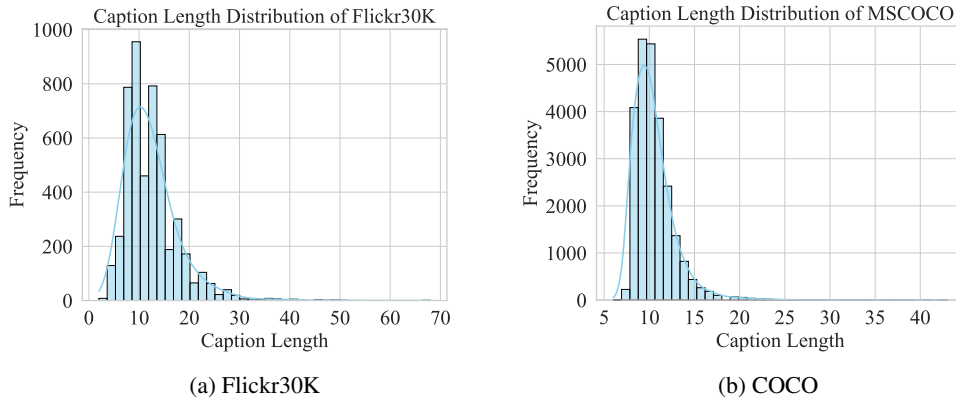


Figure 8: **Caption length distribution charts of Flickr30K and COCO.**

Table 8: **Average sentence length of generated captions of COCO-118K by InternVL2.5-38B.** Length is measured in words.

Average Length (Words)	
Left Upper	37.7
Right Upper	36.8
Left Lower	36.3
Right Lower	48.2
Summary	30.8
All	35.5

Table 9: **Image retrieval results on Flickr30K [36] and COCO [15].** We also list text retrieval results for reference. The results indicate that MLLM-based feature extractors have comparable performance to their baseline model. ReP represents Regional Prompts.

Method	Multiple Features	Image Retrieval						Text Retrieval					
		Flickr30k			COCO			Flickr30k			COCO		
		R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
E5-V [9]	✗	80.8	95.5	97.7	52.1	76.6	83.6	88.1	98.8	99.4	62.2	83.6	89.9
E5-V + Semantic Prompts	✓	80.7	95.5	97.6	<b>52.8</b>	<b>77.1</b>	<b>85.2</b>	86.8	98.1	99.0	59.6	82.4	88.8
E5-V + ReP	✓	81.0	95.7	97.8	51.8	76.3	84.4	87.9	98.7	99.4	59.9	82.0	88.6
E5-V (ft.) + Semantic Prompts	✓	<b>81.1</b>	<b>95.8</b>	97.9	51.8	76.4	84.5	91.9	<b>99.1</b>	99.7	65.1	<b>86.7</b>	<b>92.5</b>
E5-V (ft.) + ReP	✓	80.7	95.7	<b>98.0</b>	50.6	75.6	83.7	<b>93.8</b>	<b>99.1</b>	<b>99.8</b>	<b>65.5</b>	86.4	92.0

**Region-description pair annotation.** This step is performed manually by researchers ourselves. We use the online annotation tool<sup>3</sup> as the interface to draw the bounding boxes, and design the descriptions manually. In the annotation process, we try to select small objects taking up less than 10% of the image. However, there are no explicit numerical constraints in this step, as we can filter them by the area ratio to the whole image later.

**Automatic check.** We use Python scripts to automatically remove the samples where the bounding box of the target object takes up more than 20% of the image; we also automatically pick out the samples whose corresponding query descriptions are shorter than 8 words, and extend their lengths to at least 8 words.

## D More Visualization Examples in SORCE-1K

In Fig. 9, we present more samples from our benchmark for reference.

## E Re-caption Prompts for COCO-118K

The prompt we use for guiding InternVL2.5-38B [4] is as follows:

Describe the image in a structured manner, following this order: left upper, right upper, left lower, right lower, and summary. For each area, include detailed descriptions of key elements, such as objects, colors, textures, sizes, actions, or any other relevant features. Ensure each description is under 100 words, but provide enough detail to give a vivid picture. The output should be a dictionary with the following keys: ‘left upper’, ‘right upper’, ‘left lower’, ‘right lower’ and ‘summary.’

The resulting average length of generated captions is listed in Tab. 8

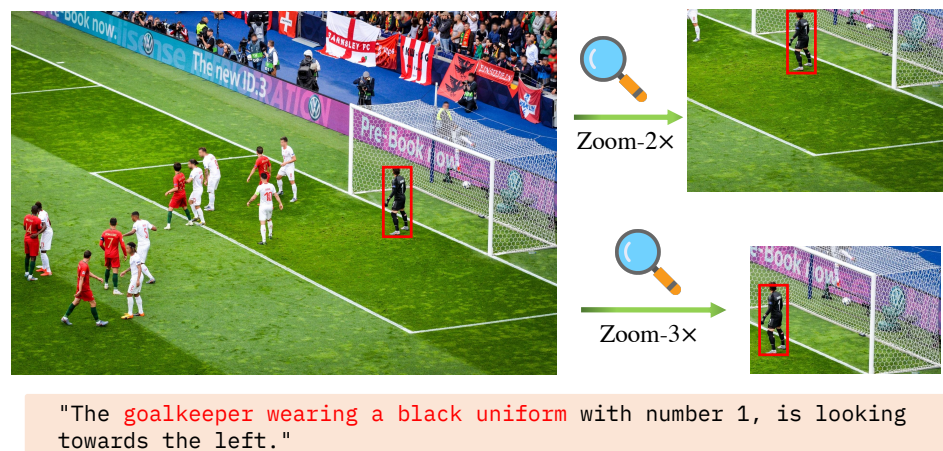
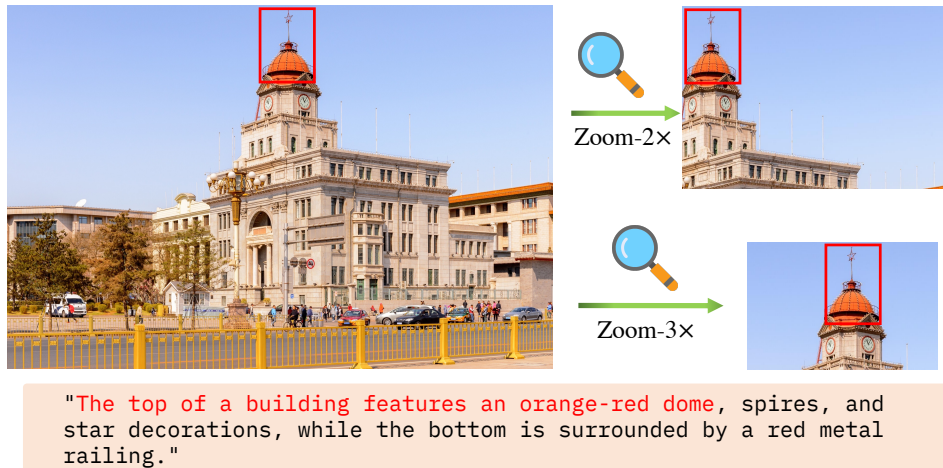
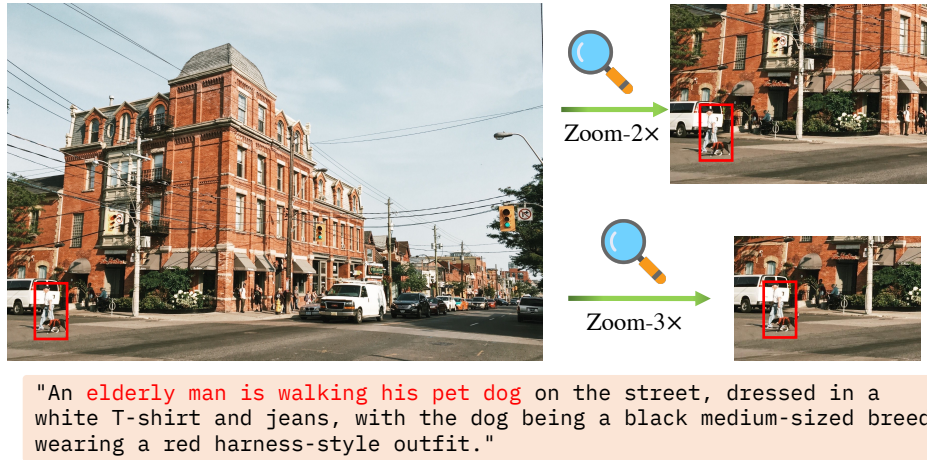


Figure 9: More visualization examples from SORCE-1K benchmark.

Table 10: **Image retrieval results on SORCE-1K.** We report R@1, R@5, and R@10. The results indicate that ReP enhances the performance of MLLM-based feature extractors, while fine-tuning further improves retrieval recall. (ft.) means fine-tuned, please refer to Section 4.2. ReP represents Regional Prompts.

Method	Multiple Features	SORCE-1K (Ours)								
		Zoom-3×			Zoom-2×			Full Res.		
		R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
E5-V [9]	✗	57.6	74.0	82.2	42.4	62.2	69.0	21.9	36.7	44.6
E5-V + Semantic Prompts	✓	57.9	74.2	81.8	42.4	61.1	68.9	21.0	36.1	45.1
E5-V + ReP	✓	62.0	80.6	86.3	54.0	73.2	80.2	27.7	45.4	53.7
E5-V (ft.) + Semantic Prompts	✓	65.8	82.8	88.1	53.5	72.4	80.2	28.0	45.7	54.6
E5-V (ft.) + ReP	✓	<b>68.0</b>	<b>85.5</b>	<b>90.9</b>	<b>56.3</b>	<b>77.1</b>	<b>83.0</b>	<b>31.5</b>	<b>50.6</b>	<b>60.0</b>

## F Exploration of Other Prompts

In addition to Regional Prompts (ReP), we also explored extracting multiple features using semantic prompts, such as “Summarize the [main component/ background/ detail] in the above image in one word:”. We append the global feature to these, resulting in a set of four features per image. For training dataset construction, we prompt the InternVL2.5-38B to generate descriptions of each image from different perspectives. We then employ contrastive fine-tuning with random choice, using the same learning hyperparameters as described in Sec. 4.2.

As shown in Tab. 9, the performance of semantic prompts remains comparable after fine-tuning. However, in Tab. 10, while semantic prompts also improve performance after fine-tuning, regional prompts consistently outperform them. We attribute this to the inherent ambiguity of semantic prompts, as it is often unclear which part of the image constitutes the "main component" or "detail." This conclusion is evident from the E5-V before finetuning, indicating that regional prompt instructions are easier for MLLMs to follow. Nevertheless, these results highlight the potential of MLLM-based text-guided feature representations for enhancing retrieval tasks.

## G Discussions and Limitations.

Although extracting multiple features by regional prompting can keep the performance of the common image retrieval, and exhibits superior performance on small objects in complex scene image retrieval. The storage requirement also multiplies. Therefore, a feature extractor with the capability to preserve every piece of visual information into one feature is still an ideal option. However, this is very challenging. Therefore, as a compromise solution, MLLM as a feature extractor is worth discussing. For instance, we could see the potential through our paper, which is that text-guided feature extraction is possible. Further, we could possibly utilize the autoregressive ability of MLLM to adaptively extract important features for the image, saving the need for extracting a fixed number of features for each image.

## H Broader Impact

In this work, we introduce the task of Small Object Retrieval in Complex Environments, namely, SORCE, the first T2IR benchmark focused on unsalient small objects. And we propose a benchmark, SORCE-1K, to evaluate the performance of the models. Furthermore, we provide an initial solution by Regional Prompts (ReP) to extract corresponding image features which can focus on different aspects of the image while keeping the global information. The SORCE task and our benchmark are beneficial for enhancing the understanding and feature extraction performance of multimodal models. Also, as a research-oriented work and the initial attempt on such a task, we trained our model on a very limited set of dataset (COCO-118K [15]), which is only for validating our idea and may not be ready for real-world applications. It may be mitigated with fine-tuning more real-world-related datasets.

<sup>3</sup><https://trexlabel.com/>

Table 11: Licenses and URLs for every dataset, benchmark, code, and pretrained models used in this paper.

Assets		License	URL
Dataset and Benchmarks	MSCOCO	CC BY 4.0	<a href="https://cocodataset.org/">https://cocodataset.org/</a>
	Flickr30K	Research purposes only	<a href="https://shannon.cs.illinois.edu/DenotationGraph/">https://shannon.cs.illinois.edu/DenotationGraph/</a>
Codes and Pretrained Models	SA-1B	Research purposes only	<a href="https://ai.meta.com/datasets/segment-anything/">https://ai.meta.com/datasets/segment-anything/</a>
	CLIP	MIT-license	<a href="https://github.com/openai/CLIP">https://github.com/openai/CLIP</a>
	EVA-02-CLIP	MIT-license	<a href="https://github.com/baaivision/EVA">https://github.com/baaivision/EVA</a>
	LLaVA-Next	Apache-2.0 license	<a href="https://github.com/LLaVA-VL/LLaVA-NeXT">https://github.com/LLaVA-VL/LLaVA-NeXT</a>
	E5-V	Research purposes only	<a href="https://github.com/kongds/E5-V">https://github.com/kongds/E5-V</a>
	InternVL2.5	MIT-license	<a href="https://github.com/OpenGVLab/InternVL/tree/main">https://github.com/OpenGVLab/InternVL/tree/main</a>

## I License of datasets and pre-trained models

All the datasets we used in the paper are commonly used datasets for academic purposes. All the licenses of the used benchmark, codes, and pretrained models are listed in Tab. 11.