

# iSEARLE: Improving Textual Inversion for Zero-Shot Composed Image Retrieval

Lorenzo Agnolucci\*, Alberto Baldrati\*, Alberto Del Bimbo, Marco Bertini

**Abstract**—Given a query consisting of a reference image and a relative caption, Composed Image Retrieval (CIR) aims to retrieve target images visually similar to the reference one while incorporating the changes specified in the relative caption. The reliance of supervised methods on labor-intensive manually labeled datasets hinders their broad applicability to CIR. In this work, we introduce a new task, Zero-Shot CIR (ZS-CIR), that addresses CIR without the need for a labeled training dataset. We propose an approach, named iSEARLE (improved zero-Shot composEd imAge Retrieval with textual invErsion), that involves mapping the visual information of the reference image into a pseudo-word token in the CLIP token embedding space and combining it with the relative caption. To foster research on ZS-CIR, we present an open-domain benchmarking dataset named CIRCO (Composed Image Retrieval on Common Objects in context), the first CIR dataset where each query is labeled with multiple ground truths and a semantic categorization. The experimental results illustrate that iSEARLE obtains state-of-the-art performance on three different CIR datasets – FashionIQ, CIRR, and the proposed CIRCO – and two additional evaluation settings, namely domain conversion and object composition. The dataset, code, and model are publicly available at <https://github.com/miccunifi/SEARLE>.

**Index Terms**—CLIP, Composed Image Retrieval, Textual Inversion, Multimodal Learning, Image Retrieval

## I. INTRODUCTION

When provided with a query consisting of a reference image and a relative caption, Composed Image Retrieval (CIR) [1], [2] seeks to retrieve target images that visually resemble the reference one while including the modifications described in the relative caption. The bi-modal structure of the query allows users to specify the desired image characteristics more precisely. It leverages the strengths of both language-based descriptions and visual features, as certain attributes are more effectively communicated through text, while others are better expressed visually. We provide some query examples in Fig. 3.

Datasets for CIR comprise triplets  $(I_r, T_r, I_t)$ , each including a reference image, a relative caption, and a target image, respectively. The creation of datasets for CIR is costly, primarily because such data is not readily accessible online, and automated generation remains a significant challenge. Consequently, researchers are compelled to undertake labor-intensive manual labeling efforts. The manual process entails the identification of reference and target image pairs and the composition of descriptive captions that outline the differences between them. Carrying out this task is both time-consuming

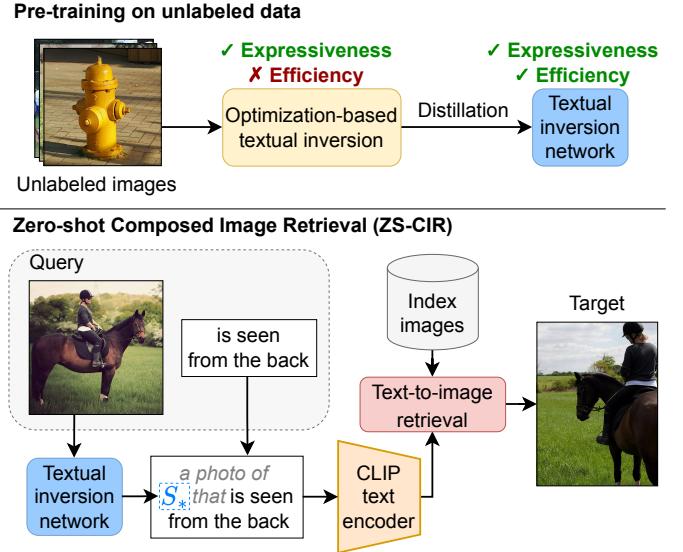


Fig. 1. Workflow of our method. *Top*: in the pre-training phase, we start by generating the pseudo-word tokens of unlabeled images with an expressive but computationally expensive optimization-based textual inversion. Then, we distill the knowledge embedded in the pseudo-word tokens into an expressive and efficient textual inversion network. *Bottom*: at inference time on ZS-CIR, we use the textual inversion network to map the reference image to pseudo-word  $S_*$  and concatenate it with the relative caption. Then, we perform text-to-image retrieval using the features extracted with the CLIP text encoder.

and resource-intensive, particularly when building extensive training datasets.

Current works addressing CIR [2]–[8] depend on supervised learning to devise methods able to combine the reference image and the relative caption effectively. For example, [5] proposes a fully supervised two-stage approach based on fine-tuning CLIP encoders and training a combiner network. Despite their promising results, the dependence of current CIR approaches on expensive manually annotated datasets constrains their scalability and applicability in domains outside those of the training datasets.

Building on the conference version of this work [9], we remove the need for costly labeled training data by introducing a new task: Zero-Shot Composed Image Retrieval (ZS-CIR). In ZS-CIR, the goal is to devise an approach capable of merging the information of the reference image and the relative caption without requiring supervised learning.

To address the challenges of ZS-CIR, we propose an approach named iSEARLE<sup>1</sup> (improved zero-Shot composEd

\* The first two authors contributed equally to this work.

L. Agnolucci, A. Baldrati, A. Del Bimbo, M. Bertini are with the Media Integration and Communication Center (MICC), University of Florence, Italy (e-mail: [name].[surname]@unifi.it).

Manuscript received May, 2024.

<sup>1</sup>John Searle is an American philosopher renowned for his work on the philosophy of language and how words denote specific objects.

imAge Retrieval with textual invErsion) based on the frozen pre-trained CLIP [10] vision-language model. Our method simplifies CIR, casting it to a standard text-to-image retrieval task by mapping the reference image into a learned pseudo-word, which is subsequently appended to the relative caption. The pseudo-word corresponds to a pseudo-word token residing in the CLIP token embedding space. We refer to this mapping process with *textual inversion*, following the terminology introduced in [11]. iSEARLE involves the pre-training of a textual inversion network – denoted as  $\phi$  – on an unlabeled image-only dataset. The pre-training process consists of two stages: an Optimization-based Textual Inversion (OTI) with a GPT-powered regularization loss aimed at generating a set of pseudo-word tokens, and the distillation of their knowledge to  $\phi$ . Upon completion of the training, the network  $\phi$  is capable of carrying out textual inversion in a single forward pass. During inference, when presented with a query  $(I_r, T_r)$ , we employ  $\phi$  to predict the pseudo-word corresponding to  $I_r$  and then concatenate it to  $T_r$ . Afterward, we exploit the CLIP common embedding space to perform text-to-image retrieval. Figure 1 shows the workflow of the proposed approach.

The majority of existing CIR datasets focus on specific domains, such as fashion [12]–[15], birds [16], or synthetic objects [1]. To the best of our knowledge, the CIRR dataset [2] stands alone in encompassing natural images within an open domain. However, CIRR suffers from two main problems. Firstly, it includes numerous false negatives, potentially leading to imprecise performance evaluations. Secondly, the queries often neglect the visual content of the reference image, rendering the task addressable through standard text-to-image techniques, as shown by the results reported in Tab. III. Additionally, existing CIR datasets provide only a single annotated ground truth image per query. To address these shortcomings and foster research on ZS-CIR, we introduce an open-domain benchmarking dataset named CIRCO<sup>2</sup> (Composed Image Retrieval on Common Objects in context). CIRCO comprises validation and test sets derived from images within the COCO dataset [17]. As a benchmarking dataset for ZS-CIR, a large training set is not needed, leading to a considerable reduction in labeling effort. To overcome the single ground truth limitation of existing CIR datasets, we propose to leverage our method to ease the annotation process of multiple ground truths. Consequently, CIRCO is the first CIR dataset with multiple annotated ground truths, enabling a more comprehensive evaluation of CIR models. In addition, contrary to existing CIR datasets, we provide a semantic categorization of the queries that allows a fine-grained semantic analysis of the results. We release only the validation set ground truths of CIRCO and host an evaluation server, enabling researchers to get performance metrics on the test set<sup>3</sup>.

The experimental results show that iSEARLE achieves state-of-the-art performance on three different CIR datasets: FashionIQ [14], CIRR [2], and the proposed CIRCO. Moreover, the experiments on two additional settings, namely domain conversion and object composition [18], prove that our model

has better generalization capabilities than competing methods.

Our contributions can be summarized as follows:

- We introduce a new task, Zero-Shot Composed Image Retrieval (ZS-CIR), to eliminate the requirement for costly labeled data for CIR;
- We propose a novel approach, named iSEARLE, that relies on a textual inversion network to address ZS-CIR by mapping images into pseudo-words. Our method comprises two phases: an optimization-based textual inversion using a GPT-powered regularization loss and the training of the textual inversion network with a distillation loss;
- We introduce CIRCO, an open-domain benchmarking dataset for ZS-CIR with multiple annotated ground truths, reduced false negatives, and a semantic categorization of the queries. We propose to leverage our model to simplify the annotation process;
- iSEARLE achieves state-of-the-art results on three different CIR datasets – FashionIQ, CIRR, and the proposed CIRCO – and two additional evaluation settings, *i.e.* domain conversion and object composition.

This work extends our conference paper [9] in several aspects: 1) we improve our method by: i) adding Gaussian noise to the text features during OTI to mitigate the issue of the *modality gap* [19]; ii) employing an additional regularization loss while training  $\phi$  to prevent the predicted pseudo-word tokens from residing in sparse regions of the CLIP token embedding space; iii) proposing a hard negative sampling strategy to help  $\phi$  in capturing fine-grained details; 2) we perform an additional annotation phase to allow a fine-grained semantic analysis on CIRCO, and we provide a more detailed study of our dataset; 3) we conduct more comprehensive experiments, including additional competitors and evaluation settings; 4) we perform a more thorough analysis of the proposed approach by studying the impact of the pre-training dataset and the effectiveness of the pseudo-word tokens in capturing visual information.

The remainder of this paper is organized as follows. Sec. II reviews related work. Sec. III details our proposed approach. Sec. IV describes the proposed CIRCO dataset. Sec. V presents experimental results and analysis. Sec. VI concludes the paper with final remarks.

## II. RELATED WORK

**Composed Image Retrieval** CIR is a branch of compositional learning, an area that has been widely explored in various vision and language tasks. These include visual question answering [20], [21], image captioning [22], [23], and image synthesis [24], [25]. Compositional learning aims to create joint embedding features that effectively integrate and express information from both the textual and visual domains.

The research on CIR spans several domains, including fashion [12]–[15], natural images [2], [16], and synthetic images [1]. The task was first introduced in [1], where the authors propose a residual gating method for composing image-text features, aiming to merge multimodal information effectively. More recently, the use of the CLIP model as a backbone for CIR has received increasing attention [3]–[6]. [4] shows the

<sup>2</sup>CIRCO is pronounced as /firkō/.

<sup>3</sup>Accessible at: <https://circo.micc.unifi.it>

effectiveness of combining out-of-the-box CLIP features with a Combiner network. Building on this, [5] introduces a task-specific fine-tuning step for CLIP encoders. Unlike the aforementioned approaches, the proposed method does not require supervision and uses unlabeled images for training, effectively learning to combine multimodal information without relying on a manually annotated CIR dataset.

**Zero-Shot Composed Image Retrieval** The Zero-Shot Composed Image Retrieval (ZS-CIR) task was introduced concurrently by Pic2Word [18] and the conference version of this work [9]. Since its introduction, several works have proposed zero-shot approaches that do not rely on costly manually annotated datasets [26]–[34].

A line of research tackles ZS-CIR by substituting the manually labeled triplets with automatically constructed ones using an LLM [26]–[28]. Specifically, [26] proposes a GPT-3-based method [35] for generating CIR triplets from an existing VQA dataset by leveraging question-answer pairs. A similar strategy is adopted by [28], which uses ChatGPT to automatically construct the triplets starting from image-caption pairs. In contrast, iSEARLE does not require any triplet-based training, as it relies only on unlabeled images. A different line of research also employs LLMs for ZS-CIR but uses them as auxiliary models at inference time rather than for automatic dataset construction [30]–[32]. For example, [31] presents a training-free approach that casts CIR to standard text-to-image retrieval by using an LLM to combine the relative caption with an automatically generated caption of the reference image. Despite the promising results, the reliance on an LLM at inference time introduces a non-negligible computational overhead when performing the retrieval.

Among the approaches addressing ZS-CIR, the most similar to our work are [18], [33], [34], as they present different methods for performing textual inversion while keeping the CLIP backbone frozen.

**Textual Inversion** In text-to-image synthesis, mapping images to a single pseudo-word is emerging as a powerful technique for generating personalized images [11], [36], [37]. [11] performs textual inversion by relying on the reconstruction loss of a latent diffusion model [24]. Additionally, [36] also fine-tunes a pre-trained text-to-image diffusion model.

Besides personalized text-to-image synthesis, textual inversion has also been applied to image retrieval tasks [18], [33], [34], [38], [39]. Specifically, PALAVRA [38] addresses personalized image retrieval by pre-training a mapping function and then optimizing the predicted pseudo-word token at inference time. Several works employ textual inversion to address ZS-CIR [18], [33], [34]. LinCIR [34] is a language-only approach for training the textual inversion network. Context-I2W [33] is based on a transformer-based textual inversion network trained on the image-caption pairs of the CC3M dataset [40]. Moreover, the textual inversion process is also dependent on the query text. This comes at the cost of requiring a double forward pass of the text encoder and a more complex network architecture than the proposed single MLP approach. The method most similar to ours is Pic2Word [18]. Pic2Word relies on a textual inversion network trained

on the 3M images of CC3M using only a cycle contrastive loss. In contrast, we train our textual inversion network on only 3% of the data and use a weighted sum of distillation and regularization losses. The distillation loss leverages the information provided by a set of pre-generated tokens obtained via optimization-based textual inversion.

**Knowledge Distillation** Knowledge distillation is a machine learning technique in which a simpler model (the student) learns to replicate the behavior of a more complex one (the teacher) by learning from its predictions [41]. This method has proven effective in various computer vision tasks, such as image classification [41]–[43], object detection [44], [45], and text-to-image synthesis [25], [46], improving model compression, computational efficiency, and accuracy. In our work, we refer to knowledge distillation as the process of transferring knowledge from a computationally expensive optimization method (teacher) to a more efficient neural network (student). Specifically, we train a textual inversion network to emulate the output of an optimization-based textual inversion using a distillation loss. From a different point of view, our lightweight network can be viewed as a surrogate model of the more resource-intensive optimization technique.

### III. PROPOSED APPROACH

**Preliminaries** CLIP (Contrastive Language-Image Pre-training) [10] is a vision and language model trained on a large-scale dataset to align images and corresponding text captions in a common embedding space. CLIP is composed of an image encoder  $\psi_I$  and a text encoder  $\psi_T$ . Given an image  $I$ , the image encoder extracts its feature representation  $x = \psi_I(I) \in \mathbb{R}^d$ , where  $d$  is the size of CLIP embedding space. For a given text caption  $T$ , a word embedding layer  $E_w$  maps each tokenized word to the token embedding space  $\mathcal{W}$ . Then, the text encoder  $\psi_T$  generates the textual feature representation  $y = \psi_T(E_w(T)) \in \mathbb{R}^d$  from the token embeddings. CLIP is trained to ensure that images and text expressing the same concepts correspond to similar feature representations within the shared embedding space.

**Overview** Starting from a frozen pre-trained CLIP model, the proposed method, named iSEARLE, is designed to generate a representation of the reference image that can be used as input to the CLIP text encoder. We achieve this goal by mapping the visual features of the image into a new token embedding within the CLIP token embedding space  $\mathcal{W}$ . We term this token embedding *pseudo-word token*, as it does not correspond to an actual word but rather serves as a representation of the image features within  $\mathcal{W}$ .

Our objective is dual. First, the pseudo-word token must accurately capture the content of the reference image. In other words, the text features related to a basic prompt comprising the pseudo-word should closely align with the corresponding image features. Second, the pseudo-word must effectively integrate and communicate with the text of the relative caption. While a single image can be mapped to multiple pseudo-word tokens, we opt for using a single one, as it proves to be sufficient to encode the information of an image effectively (see Tab. XI). Moreover, from our preliminary experiments,

### Optimization-based Textual Inversion (OTI) - Single iteration

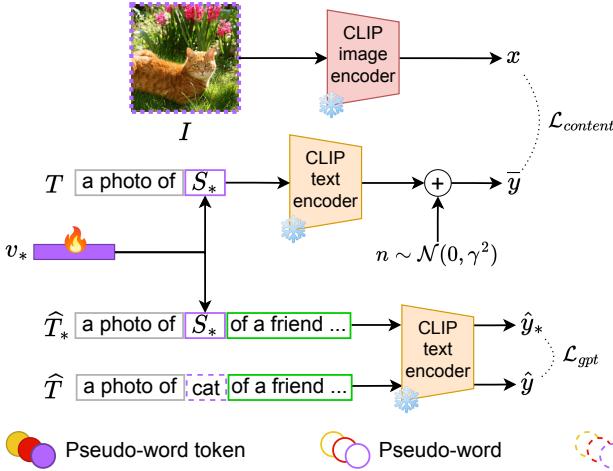


Fig. 2. Overview of our approach. *Left:* single iteration of the iterative Optimization-based Textual Inversion (OTI) method to generate a pseudo-word token  $v_*$  from an image  $I$ . We force  $v_*$  to represent the image content with a cosine loss  $\mathcal{L}_{content}$ . We add Gaussian noise to the text features before computing  $\mathcal{L}_{content}$ . We assign a concept word to  $I$  with a CLIP zero-shot classification and feed the prompt “a photo of {concept}” to GPT to continue the phrase, resulting in  $\hat{T}$ . Let  $S_*$  be the pseudo-word associated with  $v_*$ . We craft  $\hat{T}_*$  by replacing in  $\hat{T}$  the concept with  $S_*$ .  $\hat{T}$  and  $\hat{T}_*$  are then employed for a contextualized regularization with  $\mathcal{L}_{gpt}$ . *Right:* pre-training of textual inversion network  $\phi$  on unlabeled images. Given a set of pseudo-word tokens pre-generated with OTI, we distill their knowledge to  $\phi$  through a contrastive loss  $\mathcal{L}_{distil}$ . We regularize the output of  $\phi$  with the same GPT-powered loss  $\mathcal{L}_{gpt}$  employed in OTI and an additional penalty term  $\mathcal{L}_{pen}$ .  $B$  represents the number of images in a batch.

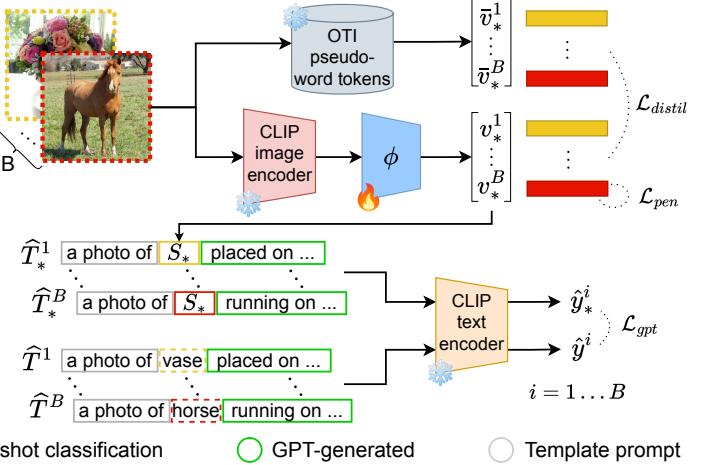
using a single pseudo-word token achieves better performance than relying on multiple ones, as also observed by [11].

iSEARLE entails pre-training a textual inversion network  $\phi$  using an unlabeled image-only dataset through a two-stage process. First, we rely on an Optimization-based Textual Inversion (OTI) method, which iteratively produces a set of pseudo-word tokens by exploiting a GPT-based regularization loss. Second, we train  $\phi$  via knowledge distillation from the pre-generated pseudo-word tokens. The  $\phi$  network outputs the pseudo-word token associated with an image in a single forward pass by taking as input its visual features, previously extracted via the CLIP image encoder.

At inference time, the input for CIR is given by a query  $(I_r, T_r)$  corresponding respectively to the reference image and the relative caption. We generate the pseudo-word token  $v_*$  associated with the reference image as  $v_* = \phi(I_r)$ . We denote by  $S_*$  the pseudo-word corresponding to the pseudo-word token  $v_*$ , *i.e.* the counterpart of  $v_*$  expressed in natural language. To effectively integrate the visual information of  $I_r$  with  $T_r$ , we build the template “a photo of  $S_*$  that {relative caption}” and extract its features via the CLIP text encoder. Note that these text features provide a multimodal representation of the reference image and its associated relative caption, as they encompass both textual and visual information. Finally, we carry out standard text-to-image retrieval within an image database using the extracted text features. We provide an overview of the workflow of our method in Fig. 1.

Fundamentally, both OTI and  $\phi$  carry out the same operation, *i.e.* mapping the visual features of an image into a pseudo-word token by means of a textual inversion. Consequently, one could directly utilize OTI at inference time without the need for  $\phi$ . However,  $\phi$  offers considerably improved efficiency compared to OTI, which needs a non-negligible amount of time to be performed (see Sec. V-A). Considering

### Pre-training of textual inversion network $\phi$



that OTI has demonstrated its effectiveness in generating fruitful pseudo-word tokens (see Sec. V), we propose to distill their knowledge into a feed-forward network. Our approach strives to maintain the powerful expressiveness of OTI while achieving a negligible inference time. From now on, we refer to our approach as iSEARLE when relying on  $\phi$  for generating the pseudo-word token and as iSEARLE-OTI when we directly utilize OTI for inference.

#### A. Optimization-based Textual Inversion (OTI)

Given an image  $I$ , we carry out textual inversion through an optimization-based approach that iteratively optimizes the pseudo-word token  $v_* \in \mathcal{W}$  for a fixed number of iterations. The left section of Fig. 2 provides an overview of OTI.

First, we randomly initialize the pseudo-word token  $v_*$  and associate the pseudo-word  $S_*$  with it. Then, we craft a template sentence  $T$ , such as “a photo of  $S_*$ ”, and process it with the CLIP text encoder  $\psi_T$ , resulting in  $y = \psi_T(T)$ . Following [38], we randomly sample  $T$  from a given set of templates. We employ the CLIP image encoder  $\psi_I$  to extract the visual features  $x = \psi_I(I)$ .

Our goal is to obtain a pseudo-word token  $v_*$  that captures the informative content of  $I$ . To this end, in our preliminary work [9], we directly minimize the discrepancy between the image and text features by leveraging the CLIP common embedding space. However, [19] shows that in vision-language models such as CLIP the features associated with text and images correspond to different regions of the joint embedding space. In other words, text and image embeddings fall into separate clusters in the feature space. This phenomenon is commonly referred to as *modality gap* [19], [47], [48]. To mitigate this issue, [48] proposes a simple and training-free strategy that involves adding Gaussian noise to the text features. Intuitively, the noise reduces the modality gap by

spreading out the text embeddings to make them overlap with the image ones.

Inspired by [48], we propose to add Gaussian noise to the text features  $y$  before minimizing their discrepancy with the image features  $x$ . Specifically, we compute  $\bar{y} = y + n$ , where  $n \sim \mathcal{N}(0, \gamma^2)$  is drawn from a Gaussian distribution with variance  $\gamma^2$ . Finally, we employ a cosine loss to maximize the similarity between the image and noisy text features:

$$\mathcal{L}_{content} = 1 - \cos(x, \bar{y}) \quad (1)$$

Therefore, differently from our preliminary work [9], we mitigate the modality gap issue before computing the loss, resulting in improved performance (see Sec. V-C). In addition, despite addressing the modality gap problem, [48] does not directly contrast noisy text features with image ones, as the authors train their model without relying on visual data. On the contrary, we show that adding Gaussian noise to the text features is effective even when they are directly compared to the image features in the loss computation.

Relying solely on  $\mathcal{L}_{content}$  is inadequate for generating a pseudo-word capable of interacting with other words of the CLIP dictionary. Indeed, similar to [38], we observe that  $\mathcal{L}_{content}$  pushes the pseudo-word token into sparse regions of CLIP token embedding space that differ from those encountered during CLIP’s training. This phenomenon, akin to effects observed in GAN inversion works [49], [50], hampers the ability of the pseudo-word token to communicate effectively with other tokens. To address this limitation, we propose a novel regularization technique that constrains the pseudo-word token to reside on the CLIP token embedding manifold, thereby enhancing its interaction capabilities. Relying on CLIP zero-shot capabilities, we carry out a zero-shot classification of the image  $I$ . To classify the images, we employ a vocabulary originating from the  $\sim 20K$  class names of the Open Images V7 dataset [51]. Specifically, we assign the  $k$  most similar distinct class names to each image, with  $k$  being a hyperparameter. We will refer to these class names as *concepts*, so, in other words, we associate each image to  $k$  different concepts. Differently from [38], we do not require the concepts as input.

After associating a set of concepts with an image, we generate a phrase using a lightweight GPT model [35]. In each iteration of the optimization process, we randomly sample one of the  $k$  concepts related to the image  $I$  and feed the prompt “a photo of {concept}” to GPT. Given that GPT is an autoregressive generative model, it manages to continue the prompt in a meaningful manner. For example, given the concept ““cat”, the GPT-generated phrase might be  $\hat{T} = “a photo of cat that is eating in front of a window”$ . In practice, since the vocabulary is known beforehand, we pre-generate all the GPT phrases for all the concepts in the vocabulary in advance. Starting from  $\hat{T}$ , we define  $\hat{T}_*$  by simply replacing the concept with the pseudo-word  $S_*$ , resulting in  $\hat{T}_* = “a photo of  $S_*$  that is eating...”$ . We extract the features of both phrases through the CLIP text encoder, obtaining  $\hat{y} = \psi_T(\hat{T})$  and  $\hat{y}_* = \psi_T(\hat{T}_*)$ . Finally, we rely on a cosine loss to maximize the similarity between the features:

$$\mathcal{L}_{gpt} = 1 - \cos(\hat{y}, \hat{y}_*) \quad (2)$$

Intuitively,  $\mathcal{L}_{gpt}$  applies a contextualized regularization that steers  $v_*$  toward the concept while considering a broader context. Indeed, compared to a generic pre-defined prompt, the GPT-generated phrases are more structured and thus similar to the relative captions used in CIR. In this way, we improve the ability of  $v_*$  to interact with human-generated text such as the relative captions.

The final loss that we use for OTI is:

$$\mathcal{L}_{OTI} = \lambda_{content} \mathcal{L}_{content} + \lambda_{OTIgpt} \mathcal{L}_{gpt} \quad (3)$$

where  $\lambda_{content}$  and  $\lambda_{OTIgpt}$  are the loss weights. Additionally, we find that applying a weight decay regularization to the pseudo-word tokens improves the effectiveness of the inversion process.

### B. Textual Inversion Network $\phi$ Pre-training

OTI proves to be effective in generating pseudo-words that not only capture the visual information of an image but also interact fruitfully with actual words. However, its iterative and optimization-based nature results in a non-negligible amount of time for its execution (see Sec. V-A). To address this issue, we propose an approach for training a textual inversion network  $\phi$  capable of predicting the pseudo-word tokens in a single forward pass by distilling knowledge from a collection of OTI pre-generated tokens. In other words,  $\phi$  serves as a more efficient surrogate model of OTI, offering a faster and computationally less demanding approximation. The right part of Fig. 2 illustrates an overview of the pre-training phase.

We aim to obtain a single model capable of inverting images from any domain without the requirement of labeled training data. Specifically, we design an MLP-based textual inversion network  $\phi$  with three linear layers, each followed by a GELU [52] activation function and a dropout layer.

Starting from an unlabeled pre-training dataset  $\mathcal{D}$ , we apply OTI to each image. Although this step is time-intensive, it is only required once, making it acceptable. This results in a collection of pseudo-word tokens, denoted as  $\bar{\mathcal{V}}_* = \{\bar{v}_*^j\}_{j=1}^N$ , where  $N$  is the total number of images in  $\mathcal{D}$ . Our goal is to distill the knowledge captured by OTI in  $\bar{\mathcal{V}}_*$  to  $\phi$ . Given an image  $I \in \mathcal{D}$ , we extract its features via the CLIP visual encoder, resulting in  $x = \psi_I(I)$ . We exploit  $\phi$  to predict the pseudo-word token  $v_* = \phi(x)$ . We minimize the distance between the predicted pseudo-word token  $v_*$  and the associated pre-generated token  $\bar{v}_* \in \bar{\mathcal{V}}_*$  while maximizing the discriminability of each token. To achieve this, we employ a symmetric contrastive loss inspired by SimCLR [38], [53]:

$$\mathcal{L}_{distil} = \frac{1}{B} \sum_{i=1}^B -\log \frac{e^{(c(\bar{v}_*^i, v_*^i)/\tau)}}{\sum_{j=1}^B e^{(c(\bar{v}_*^i, v_*^j)/\tau)} + \sum_{j \neq i} e^{(c(v_*^i, v_*^j)/\tau)}} - \log \frac{e^{(c(v_*^i, \bar{v}_*^i)/\tau)}}{\sum_{j=1}^B e^{(c(v_*^i, \bar{v}_*^j)/\tau)} + \sum_{j \neq i} e^{(c(\bar{v}_*^i, \bar{v}_*^j)/\tau)}} \quad (4)$$

Here,  $B$  is the number of images in a batch,  $c(\cdot)$  indicates the cosine similarity, and  $\tau$  is a temperature hyperparameter.

However, since the pre-training dataset  $\mathcal{D}$  comprises real-world images depicting a wide variety of subjects, a randomly sampled batch may contain significantly diverse images. In that case, it becomes trivial for the model to distinguish between the positive and negative examples, thereby reducing the effectiveness of the learning process. To avoid this issue, we propose a strategy to guarantee the inclusion of hard negative examples in every batch, which is known to improve contrastive learning performance [54], [55]. Specifically, we first perform an a priori K-Means clustering [56] of the visual features corresponding to the images comprising  $\mathcal{D}$ . Then, during training, we structure each batch such that a proportion  $\alpha$  consists of images from the same cluster, ensuring the presence of hard negative examples. The remaining fraction,  $(1 - \alpha)$ , is filled with images randomly selected from the dataset. This approach strikes a balance by introducing challenging examples into the batch while also preserving a broad diversity within the images. This strategy differs from the one we used in the conference version of this work [9], where we simply sampled the images within each batch at random. By including visually resembling examples within each batch, we encourage the model to focus on fine-grained details, improving its ability to discriminate between similar images. Consequently, as shown by the experimental results, the proposed hard negative sampling strategy improves the performance, especially on a dataset with a narrow domain such as FashionIQ [14] (see Sec. V-C for more details).

Differently from the conference version of this work [9], we employ a combination of two losses to regularize the training of  $\phi$ . First, we employ the same  $\mathcal{L}_{gpt}$  loss described in Sec. III-A. Second, inspired by [57], we propose to use an additional regularization penalty term to constrain the norm of the predicted pseudo-word tokens:

$$\mathcal{L}_{pen} = \frac{1}{B} \sum_{i=1}^B \|v_*^i\|_2^2 \quad (5)$$

This loss contributes to preventing the pseudo-word tokens generated by  $\phi$  from residing in sparse regions of the CLIP token embedding space [57]. From another point of view,  $\mathcal{L}_{pen}$  can be interpreted as a weight decay regularization term that is applied to the output of the network instead of its parameters. This aligns with the proposed OTI approach, where we apply weight decay regularization directly to the pseudo-word tokens.

The final loss for training  $\phi$  is

$$\mathcal{L}_\phi = \lambda_{distil}\mathcal{L}_{distil} + \lambda_{\phi gpt}\mathcal{L}_{gpt} + \lambda_{pen}\mathcal{L}_{pen} \quad (6)$$

with  $\lambda_{distil}$ ,  $\lambda_{\phi gpt}$ , and  $\lambda_{pen}$  representing the loss weights.

Since we do not leverage any labeled data, the training of our textual inversion network  $\phi$  is entirely unsupervised. Indeed, differently from PALAVRA [38] and Context-I2W [33], we do not require any caption and employ only raw images. Specifically, we use the unlabeled test split of the ImageNet1K [58] dataset as  $\mathcal{D}$  to pre-train  $\phi$ . It comprises 100K images without any associated labels. Compared to Pic2Word [18] and Context-I2W [33], our method uses about 3% of the data. We selected this dataset because it contains

real-world images spanning a wide variety of subjects. The experiments show that our method is robust to the choice of the  $\phi$  pre-training dataset (see Sec. V-D for more details).

#### IV. CIRCO DATASET

We recall that CIR datasets comprise triplets  $(I_r, T_r, I_t)$  composed of a reference image, relative caption, and target image (*i.e.* the ground truth), respectively.

Existing datasets often include numerous false negatives, namely images that could potentially serve as valid ground truths for a query but are not labeled as such. This issue arises because, in each query triplet, only one image is designated as the target, rendering all other images as negatives. Additionally, most datasets are confined to specialized domains, such as fashion [12]–[15], birds [16], or synthetic objects [1]. To the best of our knowledge, the CIRR dataset [2] is the sole dataset built on real-life images across an open domain. During the data collection process of CIRR, sets of 6 visually similar images are automatically generated. Subsequently, queries are devised so that both the reference and the target images belong to the same set, aiming to avoid the presence of false negatives within that particular set. However, this strategy does not guarantee the absence of false negatives throughout the entire dataset. Moreover, despite the visual similarity, the differences between images within the same set may not be easily expressible through relative captions and might necessitate absolute descriptions. This diminishes the significance of the visual information of the reference image and makes the retrieval task addressable with standard text-to-image techniques. For more details, refer to Sec. V-B.

To address these issues, we introduce an open-domain benchmarking dataset named CIRCO (Composed Image Retrieval on Common Objects in context). CIRCO is based on open-domain real-world images and is the first dataset for CIR with multiple ground truths and fine-grained semantic annotations. The whole annotation process has been carried out by the authors of this paper. To this end, we have developed a custom annotation tool that met our needs. The annotation process consists of three phases. In the first one, we build the triplets composed of a reference image, a relative caption, and a single target image. In the second one, we extend each triplet by annotating additional ground truths. In the third one, we assign semantic aspects to each query based on the relative caption. Figure 3 shows some query examples of CIRCO.

##### A. Triplets Annotation

CIRCO is based on images sourced from COCO 2017 [17] unlabeled set, which comprises 123,403 images. This dataset was suitable for our goals as it comprises open-domain real-world images that portray a wide range of subjects. In addition, we opted for the COCO unlabeled set over the training one to avoid any pre-existing model biases, as the latter is commonly used for pre-training. In the COCO labeled sets, each object in an image is categorized under one of 12 supercategories: *person*, *animal*, *sports*, *vehicle*, *food*, *accessory*, *electronic*, *kitchen*, *furniture*, *indoor*, *outdoor*, and *appliance*.

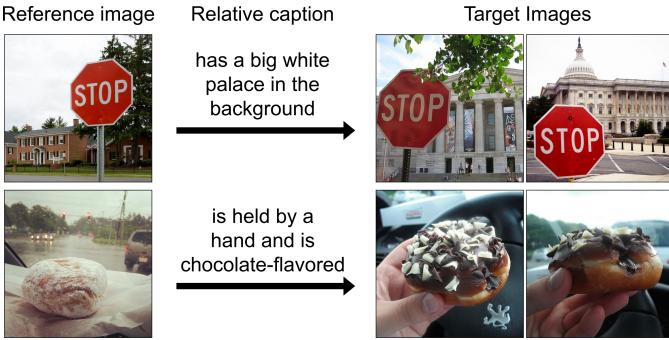


Fig. 3. Examples of CIR queries and ground truths in CIRCO.

The first step is leveraging CLIP ViT-L/14 zero-shot classification capabilities to associate every image of the unlabeled set to a supercategory. We assume that the classification is based on the predominant subject of each image. This categorization aims to get an estimation of the content of each image to be able to later create a balanced dataset. Indeed, we annotate CIRCO so that the queries comprise reference images that are evenly distributed across the supercategories. This balancing step is required to address the noticeable domain bias observed in COCO images. Indeed, certain objects, such as stop signs and fire hydrants, are over-represented.

The annotation tool selects a reference image at random and presents it alongside a gallery of 50 candidate target images. The target images must be visually similar to the reference one yet exhibit discernible disparities, as CIR requires the differences between them to be describable with a relative caption. Consequently, we choose the candidate target images based on their visual similarity to the reference image as per the CLIP features. To prevent the inclusion of near-identical images, we exclude those with a cosine similarity exceeding 0.92. The annotators are allowed to skip the current reference image if no suitable target is found in the gallery. On the contrary, when a suitable target image is available, the annotator selects it and writes the *shared concept*, which represents the common characteristics between the reference and target images. We collect the shared concept to address any potential ambiguities. Finally, the annotator crafts a relative caption from the prefix “Unlike the provided image, I want a photo of {shared concept} that”. Since our goal is to create a challenging dataset comprising truly relative captions, we ensure that they are formulated in such a way that avoids references to subjects mentioned in the shared concept. In this way, the subject of the relative caption needs to be deduced from the reference image alone.

At the end of this phase, we obtain 1020 triplets comprising a reference image, a relative caption, and a single target image.

### B. Multiple Ground Truths Annotation

For each triplet, we aim to label as ground truth all the images – besides the target one – that represent valid matches for the corresponding query. Given the starting triplet, the annotator needs to identify the ground truths from a gallery of images.

We propose to facilitate the annotation process by exploiting our approach to retrieve the images from which the ground truths are selected. In particular, we employ SEARLE to generate the pseudo-word  $S_*$  corresponding to the reference image. Then, we carry out text-to-image retrieval based on the query ‘a photo of {shared concept}  $S_*$  that {relative caption}’. During the annotation phase, we incorporate the shared concept into the query because it improves performance. Indeed, considering the single ground truth triplets obtained in Sec. IV-A, we achieve a Recall@100 of 82.15 with the shared concept and of 66.25 without it. In the gallery of images used for selecting the multiple ground truths, the annotation tool presents the top 100 retrieved images using our approach, along with the top 50 images most visually similar to the target one.

At the end of this phase, we have 4624 ground truths, of which 4097 were retrieved employing our method and 527 using the similarity with the target image. Since SEARLE achieves a Recall@100 of 82.15, by approximation, we estimate that about 82.15% of the total ground truths are present in the top 100 retrieved images. Consequently, the estimated total number of ground truths in our dataset is approximately  $4097/0.8215 \approx 4,987$ . Given that we labeled 4624 images as ground truth, we can infer that the annotated ones are  $4,624/4,987 \approx 92.7\%$  of the total. Therefore, we estimate that our annotation strategy allows us to reduce the percentage of missing ground truths in the dataset to less than 10%.

Thanks to this second annotation step, we labeled additional  $4624 - 1020 = 3604$  ground truths that would have otherwise been considered false negatives. Furthermore, this phase enables us to estimate the percentage of missing ground truths within the dataset. Notably, this estimation is unfeasible for CIR datasets featuring only a single ground truth, such as FashionIQ [14] and CIRR [2], as they lack any information about the total number of ground truths. Indeed, in these datasets, the annotation process concludes upon the completion of the triplet construction.

### C. Semantic Aspects Annotation

To allow a fine-grained semantic analysis on CIRCO, we perform a third annotation phase to assign semantic aspects to each query depending solely on its relative caption. In particular, we consider the same semantic categories as CIRR [2], such as “*direct addressing*” and “*compare & change*”.

A preliminary labeling of the semantic aspects was carried out during the multiple ground truths annotation phase in the conference version of this work [9]. The goal was to measure some raw statistics on the semantic categories of the relative captions of CIRCO. In that case, the semantic aspects were determined only by the user responsible for annotating the corresponding query. However, given the broad terms used to refer to the categories and the ambiguity of language, some inconsistencies in the classification among the annotators are inevitable. Given that we no longer aim just to compute raw statistics but rather conduct a fine-grained semantic analysis of the results on CIRCO, we require clean and reliable annotations. To this end, we introduce an

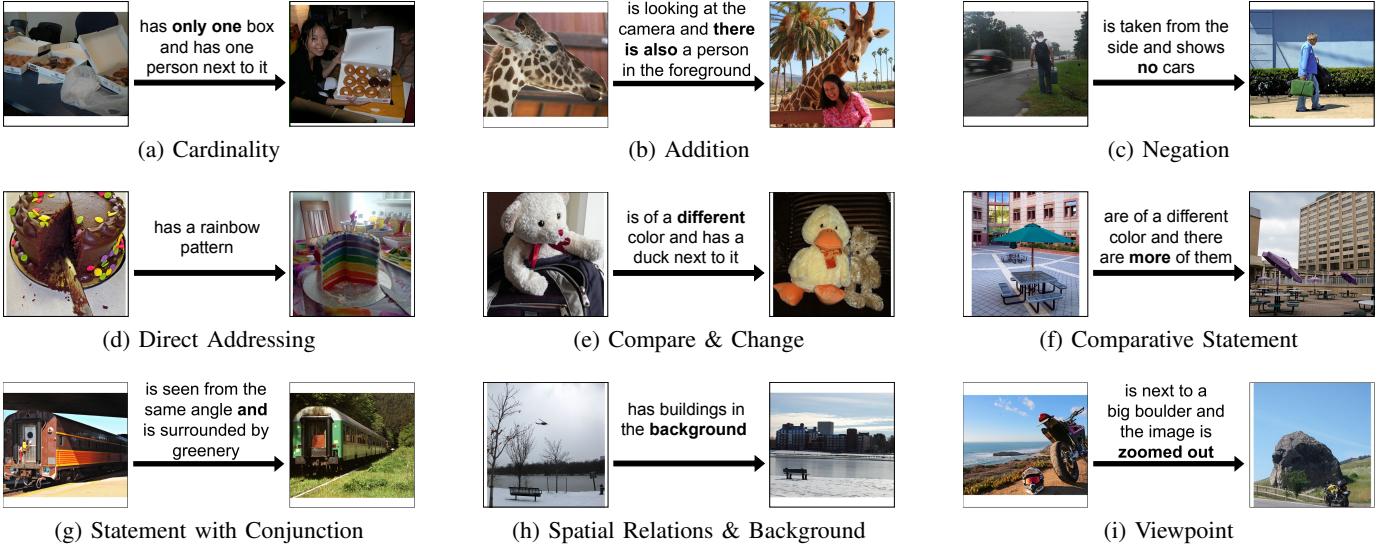


Fig. 4. Examples of queries of the proposed CIRCO dataset for different semantic aspects. For simplicity, we report only one ground truth. We highlight the keywords of each semantic aspect in bold.

additional annotation phase entirely focused on the semantic aspects related to the queries. First, we set a collection of rules on how to label each semantic category to remove possible ambiguities. In particular, a semantic aspect is assigned to a query if the relative caption: a) *cardinality*: explicitly mentions a specific number of objects in the scene, *e.g.* one, two; b) *addition*: requires the addition of an object that is not present in the reference image, *e.g.* shows also a cat; c) *negation*: requests the removal of an object present in the reference image, *e.g.* is without cars; d) *direct addressing*: explicitly requests for a specific change, *e.g.* is playing with a red ball; e) *compare & change*: requests to change something while mentioning an attribute of the reference image, *e.g.* there is a cat instead of a dog; f) *comparative statement*: includes a comparison, *e.g.* has more people; g) *statement with conjunction*: includes a conjunction proposition, *e.g.* and, or; h) *spatial relations & background*: references the background or spatial relations among objects, *e.g.* has a lake in the background; i) *viewpoint*: mentions a specific viewpoint or perspective, *e.g.* is shot from the top. Note that these categories are not mutually exclusive, *i.e.* a single relative caption can be labeled with multiple semantic aspects. We provide a query example for each semantic aspect in Fig. 4. Then, we make all the annotators label the semantic aspects of all the queries following the set of rules. Finally, we obtain the ground truth annotations by assigning a semantic aspect to a query if at least half of the annotators agree on the corresponding annotation. Therefore, differently from the conference version of this work [9], each semantic category label stems from the judgment of multiple annotators and thus has a higher reliability.

After this third annotation phase, we obtain a clean and reliable semantic categorization of the queries. As a result, CIRCO is the first CIR dataset that enables a fine-grained semantic analysis of the performance of different methods. Indeed, existing datasets [2], [14] just report raw statistics of the semantic categorization of the queries. Moreover, such

TABLE I  
ANALYSIS OF THE SEMANTIC ASPECTS COVERED BY THE RELATIVE CAPTIONS. † INDICATES RESULTS TAKEN FROM [2]. – DENOTES NO REPORTED RESULTS.

Semantic Aspect	Coverage (%)		
	CIRCO	CIRR	FashionIQ
Cardinality	16.3	29.3†	–
Addition	36.6	15.2†	15.7†
Negation	11.0	11.9†	4.0†
Direct Addressing	54.2	57.4†	49.0†
Compare & Change	37.8	31.7†	3.0†
Comparative Statement	25.7	51.7†	32.0†
Statement with Conjunction	76.2	43.7†	19.0†
Spatial Relations & Background	46.5	61.4†	–
Viewpoint	22.1	12.7†	–
Avg. Caption Length (words)	10.4	11.3†	5.3†

categorization is not publicly available, making a fine-grained semantic analysis of the performance unfeasible. In contrast, we believe that performing such an analysis of the results is crucial, as it allows us to identify the most challenging query types based on their semantic categories and thus foster focused research efforts.

#### D. Dataset Analysis

CIRCO comprises 1020 queries, randomly divided into 220 and 800 for the validation and test set, respectively. The total number of ground truths is 4624, *i.e.* 4.53 per query on average. The maximum number of ground truths annotated for a query is 21, while the modal value is 2.

The relative captions consist of an average of 10.4 words. Similar to CIRR [2], in Tab. I we report the raw statistics related to the semantic aspects associated with the relative captions. We observe that the average length of the captions

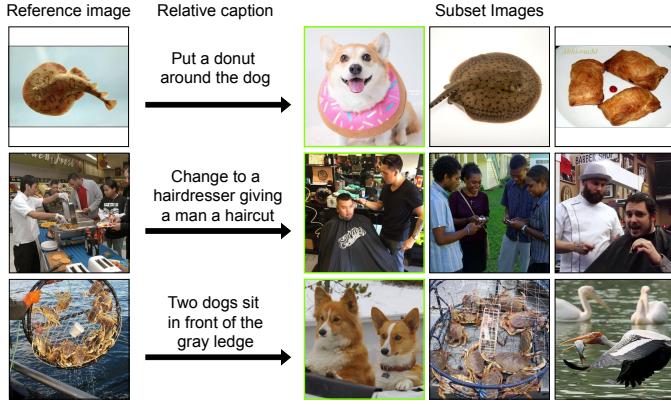


Fig. 5. Examples of queries belonging to the CIRR dataset [2]. The subset images depict very different subjects and the relative captions do not consider the reference images. We highlight the target image with a green border.

and overall coverage of the semantic categories are comparable with CIRR. However, in CIRCO approximately 75% of the annotations are composed of multiple statements, more than the ~43% of CIRR, thus revealing a higher complexity.

CIRR [2] validation and test sets comprise 4K triplets each. We remind that during the data collection process, CIRR automatically assembles subsets of 6 visually similar images based on the features of a ResNet152 [59]. Then, the queries are formulated to ensure that the reference and the target images belong to the same subset. However, despite the feature similarity, the images within these subsets often portray significantly different subjects. As a result, it becomes impossible for a human annotator to craft a relative caption, and they need to resort to an absolute description of the target image. Figure 5 shows some examples of this problem. We observe that, for instance, the annotator needs to rely on an absolute caption to describe the differences between an image depicting a crab fisherman and one with two dogs. To address this issue, we design an annotation strategy for CIRCO that lets the annotators choose the reference-target pair without constraints. As a result, we ensure that the annotators only craft captions that are truly relative, thus enhancing the quality of the dataset. To confirm this, similar to [26], we carry out an experiment to evaluate whether both the reference image and the relative caption are necessary for retrieval. Specifically, we quantitatively assess the degree of redundancy of each of the two modalities (*i.e.* image and text) by measuring the Recall@K performance of Text-to-Image (T2I) and Image-to-Image (I2I) retrieval, using respectively the relative caption and the reference image as the query. Indeed, high T2I performance implies that the relative captions are actually absolute and that the reference images are redundant. On the contrary, strong I2I results mean that the reference and target images are very similar, making the relative caption redundant. Figure 6 shows the results for varying K values. For a fair comparison with single ground truth datasets such as CIRR and FashionIQ, for CIRCO we consider only the single ground truth annotated during the first phase (Sec. IV-A). A lower curve suggests that the corresponding dataset is more difficult for a unimodal query thus indicating a lower modality redundancy.

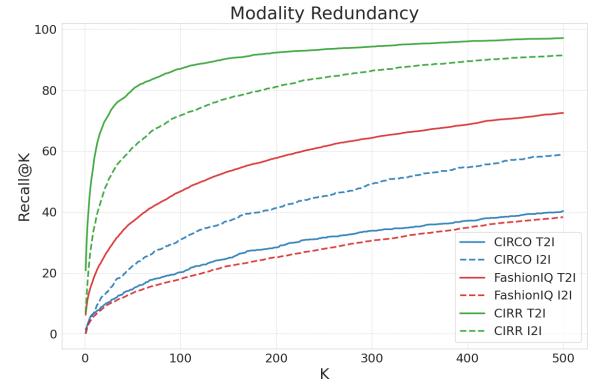


Fig. 6. Evaluation of the modality redundancy of CIRCO, CIRR, and FashionIQ validation sets. Lower values are better. T2I and I2I represent Text-to-Image and Image-to-Text retrieval, respectively.

Compared to CIRR, CIRCO demonstrates significantly lower recall metrics for both T2I and I2I, proving the quality of the proposed annotation strategy. In addition, CIRCO obtains comparable results to FashionIQ, despite encompassing a considerably broader domain.

Compared to CIRR, CIRCO comprises fewer queries, but our three-phase annotation strategy ensures higher quality, reduced false negatives, the availability of multiple ground truths, and public and reliable semantic annotations. Moreover, since its introduction in the conference version of this work [9], CIRCO has been recognized as the CIR dataset with the highest quality [30] and the cleanest annotations [31]. Finally, we employ all the 120K images of COCO as the index set, thereby providing considerably more distractors than the 2K images of the CIRR test set.

### E. Evaluation Metric

To alleviate the problem of false negatives, most works evaluate the performance on CIR datasets using Recall@K, with K set to quite large values (*e.g.* 10, 50 [14]). This makes a fine-grained analysis of the models difficult.

Thanks to the reduced false negatives and multiple ground truths of CIRCO, we can rely on a more fine-grained metric for performance evaluation, such as mean Average Precision (mAP). Indeed, mAP considers also the ranks in which the ground truths are retrieved. Specifically, we compute mAP@K, with K ranging from 5 to 50, as follows:

$$\text{mAP}@K = \frac{1}{N} \sum_{n=1}^N \frac{1}{\min(K, G_n)} \sum_{k=1}^K P@k * \text{rel}@k \quad (7)$$

where N is the number of queries,  $G_n$  is the number of ground truths of the  $n$ -th query,  $P@k$  is the precision at rank  $k$ ,  $\text{rel}@k$  is a relevance function. The relevance function is an indicator function that equals 1 if the image at rank  $k$  is labeled as a ground truth and equals 0 otherwise.

## V. EXPERIMENTAL RESULTS

We measure the performance of our method following the standard evaluation protocol [2], [3] on the three main CIR datasets: FashionIQ [14], CIRR [2] and the proposed CIRCO

TABLE II  
QUANTITATIVE RESULTS ON FASHIONIQ VALIDATION SET. BEST AND SECOND-BEST SCORES ARE HIGHLIGHTED IN BOLD AND UNDERLINED,  
RESPECTIVELY.  $\dagger$  INDICATES RESULTS FROM THE ORIGINAL PAPER.

Backbone	Method	Shirt		Dress		Toptee		Average	
		R@10	R@50	R@10	R@50	R@10	R@50	R@10	R@50
B/32	Image-only	6.92	14.23	4.46	12.19	6.32	13.77	5.90	13.37
	Text-only	19.87	34.99	15.42	35.05	20.81	40.49	18.70	36.84
	Image + Text	13.44	26.25	13.83	30.88	17.08	31.67	14.78	29.60
	Captioning	17.47	30.96	9.02	23.65	15.45	31.26	13.98	28.62
	PALAVRA [38]	21.49	37.05	17.25	35.94	20.55	38.76	19.76	37.25
	SEARLE-OTI $\dagger$ [9]	25.37	41.32	17.85	39.91	24.12	45.79	22.44	42.34
	SEARLE $\dagger$ [9]	24.44	41.61	18.54	39.51	25.70	46.46	22.89	42.53
	<b>iSEARLE-OTI</b>	<b>27.09</b>	<u>43.42</u>	<b>21.27</b>	<b>42.19</b>	<b>26.82</b>	<b>48.75</b>	<b>25.06</b>	<u>44.79</u>
L/14	<b>iSEARLE</b>	<u>25.81</u>	<b>43.52</b>	<u>20.92</u>	<b>42.19</b>	26.47	48.70	24.40	<b>44.80</b>
	Pic2Word $\dagger$ [18]	26.20	43.60	20.00	40.20	27.90	47.40	24.70	43.70
	Context-I2W $\dagger$ [33]	29.70	<u>48.60</u>	<u>23.10</u>	<u>45.30</u>	30.60	<u>52.90</u>	<u>27.80</u>	48.93
	LinCIR $\dagger$ [34]	29.10	46.81	20.92	42.44	28.81	50.18	26.28	46.49
	SEARLE-XL-OTI $\dagger$ [9]	<u>30.37</u>	47.49	21.57	44.47	30.90	51.76	27.61	47.90
	SEARLE-XL $\dagger$ [9]	26.89	45.58	20.48	43.13	29.32	49.97	25.56	46.23
	<b>iSEARLE-XL-OTI</b>	<b>31.80</b>	<b>50.20</b>	<b>24.19</b>	45.12	<b>31.72</b>	<b>53.29</b>	<b>29.24</b>	<b>49.54</b>
	<b>iSEARLE-XL</b>	28.75	47.84	22.51	<b>46.36</b>	<u>31.31</u>	52.68	27.52	<u>48.96</u>

[9]. Specifically, we use the three categories of FashionIQ validation split and the test sets of CIRR and CIRCO. Moreover, we evaluate the performance of iSEARLE on two additional settings, introduced in [18]: object composition on COCO [17] and domain conversion on ImageNet [58], [60]. In this case, we follow the evaluation protocol adopted by [18], [33].

We present two variants of our method: iSEARLE, based on CLIP ViT-B/32, and iSEARLE-XL, using CLIP ViT-L/14 as the backbone. From now on, we will refer to ViT-B/32 and ViT-L/14 as B/32 and L/14, respectively.

#### A. Implementation Details

Regarding the Optimization-based Textual Inversion (OTI), we perform 500 iterations with a learning rate of  $2e-2$ . We set the loss weights  $\lambda_{content}$  and  $\lambda_{OTIgpt}$  in Eq. (3) to 1 and 0.5, respectively. We set the standard deviation  $\gamma$  of the Gaussian noise to 0.64 and 0.16 respectively for iSEARLE-OTI and iSEARLE-XL-OTI. For the textual inversion network  $\phi$ , we train for 115 epochs, with a learning rate of  $1e-4$  and a batch size of 256. We set the loss weights  $\lambda_{distil}$  and  $\lambda_{\phi gpt}$  in Eq. (6) to 1 and 0.75, respectively. The loss weight  $\lambda_{pen}$  is equal to  $3e-3$  and  $1e-2$  for the B/32 and L/14 backbones, respectively. The temperature  $\tau$  in Eq. (4) is set to 0.25. The parameter  $\alpha$  of the hard negative sampling strategy (Sec. III-B) is set to 0.5. For both OTI and  $\phi$ , we employ the AdamW optimizer [61] with weight decay 0.01. We use an exponential moving average of 0.99 and 0.999 decay for OTI and  $\phi$ , respectively. During OTI, we set the number of concept words  $k$  associated with each image to 15, while during the training of  $\phi$  to 150. We tune each hyperparameter individually with a grid search on the CIRR validation set.

Using a single NVIDIA A100 40GB GPU, iSEARLE-XL-OTI takes  $\sim$ 35 seconds for a single image and  $\sim$ 1.1 seconds per image with batch size 256. The training of  $\phi$  for iSEARLE-XL takes 12 hours in total on a single A100 GPU. Throughout all the experiments, we adopt the pre-processing technique

introduced in [5]. For retrieval, we normalize both the query and index set features to have a unit  $L_2$ -norm.

To generate the phrases used for the regularization with  $\mathcal{L}_{gpt}$ , we exploit the GPT-Neo-2.7B model with 2.7 billion parameters developed by EleutherAI. For each of the 20,932 class names of the Open Images V7 dataset [51], we generate 256 phrases a priori with a temperature of 0.5 and a maximum length constraint of 35 tokens. The whole process requires about 12 hours to execute on a single A100 GPU. Since this operation only needs to be performed once, the time requirements are manageable.

We use the same set of templates during both training and inference. Since the FashionIQ dataset provides two relative captions per triplet, at inference time, we concatenate them using the conjunction “and”. To ensure our method remains invariant to the concatenation order, we use both possible concatenation orders and average the resulting features.

#### B. Quantitative Results

We provide the results of both iSEARLE and iSEARLE-OTI. We compare our method with several zero-shot baselines: 1) *Text-only*: we compute the similarity using only the CLIP features of the relative caption; 2) *Image-only*: we retrieve the most similar images to the reference one; 3) *Image + Text*: we sum together the CLIP features of the reference image and the relative caption; 4) *Captioning*: we substitute the pseudo-word token with the caption of the reference image obtained via a pre-trained captioning model [62]<sup>4</sup>. In addition, we compare the proposed approach with the previous version of our method [9] and state-of-the-art ZS-CIR methods: Pic2Word [18], Context-I2W [33], and LinCIR [34]. For a fair comparison, we only consider competing methods that rely on the CLIP model without fine-tuning its weights and that, at inference time, do not require any additional pre-trained models, such as LLMs.

<sup>4</sup><https://huggingface.co/laion/CoCa-ViT-B-32-laion2B-s13B-b90k>

TABLE III  
QUANTITATIVE RESULTS ON CIRR TEST SET.  $\dagger$  INDICATES RESULTS FROM THE ORIGINAL PAPER. – DENOTES RESULTS NOT REPORTED IN THE ORIGINAL PAPER.

Backbone	Method	Recall@K			
		K = 1	K = 5	K = 10	K = 50
B/32	Image-only	6.89	22.99	33.68	59.23
	Text-only	21.81	45.22	57.42	81.01
	Image + Text	11.71	35.06	48.94	77.49
	Captioning	12.46	35.04	47.71	77.35
	PALAVRA [38]	16.62	43.49	58.51	83.95
	SEARLE-OTI $\dagger$ [9]	24.27	53.25	66.10	88.84
	SEARLE $\dagger$ [9]	24.00	53.42	66.82	89.78
	<b>iSEARLE-OTI</b>	<b>26.19</b>	<b>55.18</b>	<b>68.55</b>	90.65
	<b>iSEARLE</b>	<b>25.23</b>	<b>55.69</b>	<b>68.05</b>	<b>90.82</b>
	Pic2Word $\dagger$ [18]	23.90	51.70	65.30	87.80
L/14	Context-I2W $\dagger$ [33]	<b>25.60</b>	<b>55.10</b>	<b>68.50</b>	<b>89.80</b>
	LinCIR $\dagger$ [34]	25.04	53.25	66.68	–
	SEARLE-XL-OTI $\dagger$ [9]	24.87	52.31	66.29	88.58
	SEARLE-XL $\dagger$ [9]	24.24	52.48	66.29	88.84
	<b>iSEARLE-XL-OTI</b>	<b>25.40</b>	<b>54.05</b>	<b>67.47</b>	<b>88.92</b>
	<b>iSEARLE-XL</b>	25.28	54.00	66.72	88.80

**FashionIQ** We report the results for FashionIQ in Tab. II. Considering the B/32 backbone, iSEARLE obtains comparable performance to iSEARLE-OTI, thereby preserving effectiveness while offering a notable efficiency improvement. Both versions of our approach outperform the baselines, including the preliminary version of this work [9]. Notably, the improvement over Captioning highlights that the pseudo-word token encapsulates more information than the actual words forming the generated caption. Regarding the L/14 backbone, we notice that, despite using only 3% of the training data, iSEARLE-XL achieves a considerable performance improvement over Pic2Word and comparable results with Context-I2W. In addition, we recall that Context-I2W requires a double forward pass of the text encoder and employs a transformer-based architecture significantly more complex than our MLP-based one. We observe a performance gap compared to iSEARLE-XL-OTI. We suppose that this discrepancy may stem from the very narrow domain of FashionIQ, which differs considerably from the natural images of the pre-training dataset we employ for training  $\phi$ . To support this hypothesis, we trained a version of iSEARLE-XL using the FashionIQ training set as the pre-training dataset, yielding an average Recall@10 and Recall@50 of 29.07 and 49.67, respectively, on the validation set. These results closely align with those of iSEARLE-XL-OTI, confirming our theory. We provide more details on the impact of the  $\phi$  pre-training dataset in Sec. V-D.

**CIRR** Table III shows the results for the CIRR test set. We notice that the Text-only baseline outperforms Image-only and Image+Text. These results reveal a major issue with CIRR: the relative captions are often not truly relative in practice. In particular, as observed also in [18], we notice that the reference image may not provide useful information for retrieval and may even have a detrimental effect.

We observe that, for both backbones, the results achieved by our method with OTI and  $\phi$  are comparable, thereby proving the effectiveness of the proposed distillation process.

TABLE IV  
QUANTITATIVE RESULTS ON CIRCO TEST SET.  $\dagger$  INDICATES RESULTS FROM THE ORIGINAL PAPER.

Backbone	Method	mAP@K			
		K = 5	K = 10	K = 25	K = 50
B/32	Image-only	1.34	1.60	2.12	2.41
	Text-only	2.56	2.67	2.98	3.18
	Image + Text	2.65	3.25	4.14	4.54
	Captioning	5.48	5.77	6.44	6.85
	PALAVRA [38]	4.61	5.32	6.33	6.80
	SEARLE-OTI $\dagger$ [9]	7.14	7.83	8.99	9.60
	SEARLE $\dagger$ [9]	9.35	9.94	11.13	11.84
	<b>iSEARLE-OTI</b>	<b>10.31</b>	<b>10.94</b>	<b>12.27</b>	<b>13.01</b>
L/14	<b>iSEARLE</b>	<b>10.58</b>	<b>11.24</b>	<b>12.51</b>	<b>13.26</b>
	Pic2Word [18]	8.72	9.51	10.64	11.29
	LinCIR $\dagger$ [34]	<b>12.59</b>	<b>13.58</b>	<b>15.00</b>	<b>15.85</b>
	SEARLE-XL-OTI $\dagger$ [9]	10.18	11.03	12.72	13.67
	SEARLE-XL $\dagger$ [9]	11.68	12.73	14.33	15.12
	<b>iSEARLE-XL-OTI</b>	11.31	12.67	14.46	15.34
	<b>iSEARLE-XL</b>	<b>12.50</b>	<b>13.61</b>	<b>15.36</b>	<b>16.25</b>

Interestingly, there is no performance gap between the B/32 and L/14 versions, and actually, the B/32 even outperforms the L/14 in most cases. When compared with the conference version of this work [9], our method obtains better results for both the OTI and  $\phi$  versions, highlighting the importance of the improvements introduced in this work. Regarding the L/14 backbone, Context-I2W [33] achieves the best performance. However, such a method is trained on the CC3M [40] dataset, which is more than 30 times larger than our pre-training dataset. When considering half of the training images, the authors of Context-I2W report significantly lower results, with a  $R@1$  and  $R@5$  of 24.80 and 53.60, respectively. Therefore, despite using 15 times fewer data and no captions, iSEARLE-XL still obtains comparable performance, with a  $R@1$  and  $R@5$  of 25.28 and 54.00, respectively.

**CIRCO** In Table IV, we report the results for the CIRCO test set. Firstly, we observe that, in contrast to FashionIQ and CIRR, Image+Text outperforms Image-only and Text-only. This result indicates that CIRCO contains queries where both the reference image and the relative caption are equally crucial for retrieving the target images. Secondly, iSEARLE achieves a considerable improvement over all the baselines and even outperforms Pic2Word, despite a smaller backbone. Considering the L/14 backbone, iSEARLE-XL would achieve the best results. However, we recall that the conference version of our method SEARLE-XL [9] was used to ease the annotation process of CIRCO. Therefore, since the core of the approach proposed in this work is the same, it is likely that the results could exhibit some sort of bias. Still, we report them for completeness.

Table V shows the  $mAP@10$  results on the CIRCO test set for each semantic category. We observe that some semantic aspects, such as *viewpoint* and *negation*, pose a significant challenge for all the reported methods. We suppose this outcome is due to the use of CLIP, which struggles to comprehend specific language constructs, such as negations and compositional relationships between objects and attributes [63], [64].

TABLE V  
QUANTITATIVE RESULTS ON CIRCO TEST SET FOR EACH SEMANTIC CATEGORY.

Semantic Aspect	ViT-B/32				ViT-L/14				
	PALAVRA	SEARLE	iSEARLE-OTI	iSEARLE	Pic2Word	LinCIR	SEARLE-XL	iSEARLE-XL-OTI	iSEARLE-XL
Cardinality	3.38	7.94	<u>8.29</u>	<b>9.50</b>	9.20	<u>11.80</u>	10.25	10.30	<b>11.59</b>
Addition	5.66	10.55	<b>11.95</b>	<u>11.64</u>	10.04	<b>14.66</b>	13.82	13.26	<b>14.66</b>
Negation	5.96	6.72	<b>8.51</b>	<u>7.48</u>	6.97	<b>9.91</b>	8.84	<u>9.42</u>	8.82
Direct Addressing	5.55	11.53	<u>12.41</u>	<b>12.94</b>	10.59	15.18	14.84	<u>15.29</u>	<b>16.02</b>
Compare & Change	4.30	<b>8.09</b>	8.02	<u>8.08</u>	7.48	<b>9.52</b>	<u>9.48</u>	8.30	9.42
Comparative Statement	5.82	8.38	<b>10.16</b>	<u>9.81</u>	8.47	<b>12.10</b>	11.19	10.27	<u>11.60</u>
Statement w/ Conjunction	5.25	9.35	<u>10.46</u>	<b>10.54</b>	8.94	<u>13.20</u>	12.73	12.76	<b>13.47</b>
Spatial Rel. & Background	5.89	11.30	<u>11.74</u>	<b>12.48</b>	9.97	<u>15.24</u>	14.18	14.25	<b>15.75</b>
Viewpoint	4.07	7.42	<b>7.98</b>	<u>7.45</u>	4.52	7.86	<b>8.51</b>	8.14	<u>8.25</u>

On the other hand, all the considered methods seem to handle semantic categories such as *addition* and *direct addressing* more effectively. We argue that this result stems from the fact that the relative captions corresponding to these semantic aspects have a structure similar to that of the absolute captions employed for pre-training CLIP. Regarding the comparison between different approaches, the considerations we made for Tab. IV still apply, with the proposed method achieving state-of-the-art performance across different semantic aspects.

Thanks to the semantic annotation phase introduced in this work, CIRCO allows such a fine-grained analysis of the results. Consequently, it is possible to discern the intricate complexities inherent in different query types, thereby guiding targeted research efforts toward tackling these challenges.

**Domain Conversion** Table VI illustrates the results for the domain conversion task. Following [18], the query images are sourced from 200 classes of ImageNet [58] validation set, while the target ones belong to ImageNet-R [60]. We recall that we relied on the ImageNet unlabeled test set as the pre-training dataset of our method, so there is no overlap with the evaluation images of the domain conversion task. The purpose of this experiment is to study how our model can convert the domain of a query image by using the prompt “*{domain}* of *S<sub>\*</sub>*”, where *{domain}* is a word that indicates the domain, *e.g.* *toy* or *origami*. We consider the retrieved image correct if its class is the same as that of the query image and its domain matches the one specified by the prompt. For instance, given the domain “*toy*” and a query image containing a real-world shark, the goal is to retrieve images depicting a shark toy. The results show that, both for the B/32 and L/14 backbones, our method outperforms all the baselines.

**Object Composition** We provide the results for the object composition task on the COCO validation set [17] in Tab. VII. This task aims to retrieve an image comprising an object specified with a single query image and other objects described with text. The object composition task closely resembles the personalized retrieval one [38], differing mainly in that the latter involves queries composed of multiple images depicting the same object instance. Following [18], we use the prompt “a photo of *S<sub>\*</sub>*, {*obj<sub>1</sub>*} and {*obj<sub>2</sub>*}, ..., and {*obj<sub>n</sub>*}”, where {*obj<sub>i</sub>*} are text descriptions of objects, *e.g.* *mouse*, *laptop* or *kite*. For both the B/32 and L/14 backbones, we notice that our approach achieves state-of-the-art results. Considering the

L/14 backbone, we observe a performance gap between the OTI and  $\phi$  variants of the proposed method. We suppose this result is due to the similarity between the object composition and personalized retrieval tasks, where it has been shown that an optimization-based method is more suitable and achieves better performance [38]. The results obtained by PALAVRA [38] on the object composition task confirm our hypothesis, as it achieves significantly better relative performance than in composed image retrieval, *e.g.* on CIRCO. In addition, we observe that LinCIR [34] obtains considerably worse performance than iSEARLE-XL. We suppose this outcome is due to their language-only training strategy, which makes their model struggle to capture fine-grained visual details.

**Discussion** Our experiments show that the proposed approach consistently achieves commendable performance across different datasets, highlighting its robustness. Compared to the baselines, our method has better generalization and adaptability capabilities to tasks beyond standard composed image retrieval, such as the domain conversion and object composition. Moreover, iSEARLE and iSEARLE-OTI obtain comparable results in most scenarios. This confirms the effectiveness of the distillation process, which offers a significant efficiency improvement without sacrificing the performance. Finally, the fine-grained evaluation on CIRCO reveals the intrinsic limitations of current ZS-CIR methods, which struggle when dealing with queries involving complex semantic aspects such as negation and viewpoint changes.

### C. Ablation Studies

We conduct extensive ablation studies to measure the individual contribution of each component of our approach. To avoid potential interferences, we evaluate the two main stages of the proposed method separately. In particular, we assess the performance of the textual inversion network  $\phi$  while keeping fixed the collection of OTI pre-generated tokens obtained as detailed in Sec. III-A. As  $\phi$  distills the knowledge of the OTI pre-generated tokens, we assume that the more informative they are (*i.e.* the better OTI performs), the better the performance achieved by  $\phi$  will be. We rely on the CIRR and FashionIQ validation sets to conduct the ablation studies and report the results for the main evaluation metrics. Specifically, for FashionIQ we report the average scores. We focus solely on the B/32 version of our method for simplicity.

TABLE VI  
QUANTITATIVE RESULTS FOR THE DOMAIN CONVERSION TASK. THE QUERY IMAGES ARE FROM THE IMAGENET VALIDATION SET, WHILE THE TARGET ONES BELONG TO IMAGENET-R.  $\dagger$  INDICATES RESULTS FROM THE ORIGINAL PAPER.

Backbone	Method	Cartoon		Origami		Toy		Sculpture		Average	
		R@10	R@50								
B/32	Image-only	0.22	3.09	0.43	2.38	0.55	4.53	0.47	3.80	0.42	3.45
	Text-only	0.16	1.14	1.17	5.29	0.31	1.14	0.31	1.83	0.49	2.35
	Image + Text	1.50	8.81	1.66	7.05	1.00	7.43	1.25	7.68	1.35	7.74
	Captioning	6.75	18.58	9.60	<u>21.22</u>	5.95	17.23	7.18	18.16	7.37	18.80
	PALAVRA [38]	2.56	10.81	3.29	11.39	1.48	9.44	2.89	12.50	2.56	11.04
	SEARLE-OTI [9]	7.10	18.97	8.91	19.88	5.37	17.08	6.81	18.00	7.05	18.48
	SEARLE [9]	6.12	20.24	7.91	20.03	3.10	16.18	4.56	17.35	5.42	18.45
L/14	<b>iSEARLE-OTI</b>	<b>9.49</b>	<b>24.10</b>	<b>9.93</b>	<b>21.27</b>	<b>6.96</b>	<b>21.43</b>	<b>9.21</b>	<b>22.50</b>	<b>8.90</b>	<b>22.33</b>
	<b>iSEARLE</b>	<b>10.02</b>	<b>25.01</b>	<b>9.77</b>	21.13	<b>7.07</b>	<b>22.97</b>	<b>9.16</b>	<b>22.93</b>	<b>9.01</b>	<b>23.01</b>
	Pic2Word $\dagger$ [18]	8.00	21.90	13.50	25.60	8.70	21.60	10.00	23.80	10.10	23.20
	Context-I2W $\dagger$ [33]	10.20	26.10	17.50	28.70	<u>11.60</u>	27.40	12.10	28.20	12.90	27.60
	LinCIR [34]	<u>11.34</u>	28.96	17.15	30.31	<b>13.40</b>	<u>30.30</u>	13.19	28.43	<u>13.77</u>	29.50
	SEARLE-XL-OTI [9]	9.85	24.97	18.81	30.55	10.19	27.26	12.75	28.94	12.90	27.93
	SEARLE-XL [9]	9.67	29.94	19.48	34.12	7.45	26.75	11.57	33.31	12.04	31.03
L/14	<b>iSEARLE-XL-OTI</b>	10.48	<u>29.76</u>	<u>20.01</u>	<u>33.36</u>	9.68	30.28	<u>13.39</u>	<u>33.72</u>	13.39	<u>31.78</u>
	<b>iSEARLE-XL</b>	<b>12.84</b>	<b>31.67</b>	<b>22.17</b>	<b>34.51</b>	11.20	<b>31.87</b>	<b>15.88</b>	<b>35.49</b>	<b>15.52</b>	<b>33.39</b>

TABLE VII  
QUANTITATIVE RESULTS FOR THE OBJECT COMPOSITION TASK ON COCO VALIDATION SET.  $\dagger$  INDICATES RESULTS FROM THE ORIGINAL PAPER.

Backbone	Method	R@1	R@5	R@10
B/32	Image-only	7.30	13.64	17.19
	Text-only	4.93	14.06	21.51
	Image + Text	8.96	18.24	23.92
	Captioning	5.96	15.84	22.92
	PALAVRA [38]	<u>12.94</u>	25.66	32.40
	SEARLE-OTI [9]	<b>13.03</b>	26.00	34.27
	SEARLE [9]	10.91	24.58	33.28
L/14	<b>iSEARLE-OTI</b>	11.96	<b>26.63</b>	<u>35.43</u>
	<b>iSEARLE</b>	11.75	<u>26.40</u>	<b>35.87</b>
	Pic2Word $\dagger$ [18]	11.50	24.80	33.40
	Context-I2W $\dagger$ [33]	13.50	28.50	38.10
	LinCIR [34]	10.93	24.83	34.48
	SEARLE-XL-OTI [9]	<u>17.04</u>	<u>31.43</u>	<u>40.81</u>
	SEARLE-XL [9]	14.21	29.02	37.71
<b>iSEARLE-XL-OTI</b>	<b>17.54</b>	<b>32.55</b>	<b>41.22</b>	
	<b>iSEARLE-XL</b>	15.01	30.05	38.76

**Optimization-based textual inversion (OTI)** We ablate each of the components of the optimization process: 1) *w/o GPT reg*: we regularize with a prompt containing only the concept word, without the GPT-generated suffix; 2) *random reg*: we additionally substitute the concept word with a random word; 3) *w/o reg*: we completely remove the regularization loss; 4) *w/o noise*: we do not add Gaussian noise to the text features, *i.e.* we set  $\gamma$  to 0; 5) *L<sub>2</sub> loss*: we substitute the cosine loss in Eqs. (1) and (2) with an *L<sub>2</sub>*-based one.

The upper part of Tab. VIII shows the results. As a different loss leads to a different speed of convergence, we use a tailored number of optimization iterations for each ablation experiment and report the best performance. First, we find that

regularization plays a crucial role in ensuring that the pseudo-word tokens reside in the CLIP token embedding manifold and can effectively interact with the CLIP vocabulary tokens. In particular, we argue that our GPT-based regularization loss allows the pseudo-word tokens to interact with text resembling human-written language, thereby improving their communication with the relative captions and ultimately enhancing retrieval performance. This effect is particularly pronounced in CIRR, where relative captions tend to be more elaborate and have a more diverse vocabulary. Then, we observe that using a cosine loss obtains better performance than an *L<sub>2</sub>* one. We suppose that this outcome stems from the CLIP training strategy, which uses a cosine similarity-based loss. Finally, we notice that adding Gaussian noise to the text features in the  $\mathcal{L}_{content}$  loss computation improves the performance. This result shows that our strategy for mitigating the effect of the modality gap is fruitful.

**Textual inversion network  $\phi$**  We ablate the losses we use during the pre-training of  $\phi$ : 1) *cos distil*: we use a cosine distillation loss instead of a contrastive one; 2) *w/o distil*: we replace  $\mathcal{L}_{distil}$  with the cycle contrastive loss introduced by [38], which directly considers the image and text features; 3) *w/o reg*: we remove the  $\mathcal{L}_{gpt}$  regularization loss; 4) *w/o  $\mathcal{L}_{pen}$* : we remove the regularization penalty term on the predicted pseudo-word tokens, *i.e.* we set  $\lambda_{pen}$  to 0; 5) *w/o HNSS*: we compose batches randomly instead of using the hard negative sampling strategy described in Sec. III-B, *i.e.* we set  $\alpha$  to 0.

We report the results in the lower section of Tab. VIII. The contrastive version of the distillation loss achieves better performance than the cosine one. Compared to the cycle contrastive loss, our distillation-based loss proves to be significantly more effective, highlighting how learning from OTI pre-generated tokens is more fruitful than learning from raw images. Moreover, although the pre-generated pseudo-word tokens are already regularized, we observe that our GPT-based

TABLE VIII  
ABLATION STUDIES ON CIRR AND FASHIONIQ VALIDATION SETS. FOR FASHIONIQ, WE CONSIDER THE AVERAGE RECALL.

Abl.	Method	FashionIQ		CIRR		
		R@10	R@50	R@1	R@5	R@10
OTI	w/o GPT reg	21.63	41.40	21.04	50.99	64.21
	random reg	21.53	40.01	21.21	48.43	62.97
	w/o reg	19.29	37.10	17.43	46.85	60.65
	w/o noise	23.62	43.80	24.75	55.39	69.28
	$L_2$ loss	24.40	44.71	25.23	56.80	<b>70.46</b>
	<b>iSEARLE-OTI</b>	<b>25.06</b>	<b>44.79</b>	<b>25.57</b>	<b>57.11</b>	<b>70.46</b>
$\phi$	cos distil	22.46	42.26	24.80	54.36	68.00
	w/o distil	19.64	38.54	21.93	49.94	63.55
	w/o reg	24.33	<b>45.08</b>	24.63	56.20	69.22
	w/o $\mathcal{L}_{pen}$	24.19	44.70	25.66	57.14	70.22
	w/o HNSS	23.63	44.46	25.70	<b>57.45</b>	70.24
	<b>iSEARLE</b>	<b>24.40</b>	<b>44.80</b>	<b>25.74</b>	<b>57.35</b>	<b>70.32</b>

TABLE IX  
ABLATION STUDIES ON THE CHOICE OF HYPERPARAMETERS ON CIRR AND FASHIONIQ VALIDATION SETS. FOR FASHIONIQ, WE CONSIDER THE AVERAGE RECALL.

Abl.	Value	FashionIQ		CIRR		
		R@10	R@50	R@1	R@5	R@10
OTI	$\gamma = 0.32$	24.76	<b>44.73</b>	25.16	<b>56.92</b>	69.98
	$\gamma = 0.64$	<b>25.06</b>	<b>44.79</b>	<b>25.57</b>	<b>57.11</b>	<b>70.46</b>
	$\gamma = 0.96$	24.15	44.57	<b>25.60</b>	56.54	<b>70.46</b>
$\phi$	$\alpha = 0.25$	23.76	44.49	<b>25.71</b>	<b>57.42</b>	70.28
	$\alpha = 0.5$	<b>24.40</b>	<b>44.80</b>	<b>25.74</b>	<b>57.35</b>	<b>70.32</b>
	$\alpha = 1$	<b>24.52</b>	<b>44.95</b>	25.14	56.83	70.06

regularization loss is still beneficial for training  $\phi$ , especially on the CIRR dataset. Regarding  $\mathcal{L}_{pen}$ , we find that introducing an additional regularization penalty term to constrain the predicted pseudo-word tokens is effective, as it helps in making them reside in the CLIP manifold [57]. Finally, we observe that by removing the proposed hard negative sampling strategy we achieve comparable performance on CIRR but worse on FashionIQ. This outcome is due to the narrow domain of FashionIQ, as images are closely related and subtle details are crucial. This confirms that our hard negative sampling strategy helps the model in capturing fine-grained details.

**Hyperparameters** We study the effect of the main hyperparameters of our method and report the results in Tab. IX. Regarding the standard deviation  $\gamma$  of the noise, we observe a sweet spot ( $\gamma = 0.64$ ) that balances mitigating the modality gap and preserving the informative content of the text features. For the ratio  $\alpha$  of hard negative samples per batch, composing batches entirely of hard negatives ( $\alpha = 1$ ) leads the model to focus on fine-grained details, improving performance on fine-grained datasets like FashionIQ. However, this comes at the cost of reduced image diversity within each batch, resulting in worse performance on broader-domain datasets such as CIRR. A moderate value ( $\alpha = 0.5$ ) achieves a better trade-off by incorporating challenging examples while maintaining a diverse range of image content.

TABLE X  
EVALUATION OF THE EFFECT OF THE  $\phi$  PRE-TRAINING DATASET ON THE RESPECTIVE EVALUATION SPLIT OF FASHIONIQ, CIRR, AND CIRCO. FOR FASHIONIQ, WE CONSIDER THE AVERAGE RECALL.

Method	FashionIQ		CIRR			CIRCO
	R@10	R@50	R@1	R@5	R@10	mAP@10
iSEARLE-FIQ	<b>25.61</b>	<b>45.57</b>	25.01	55.06	67.90	10.77
iSEARLE-CIRR	23.75	44.07	25.21	54.24	<b>68.15</b>	10.49
iSEARLE-NABirds	23.27	42.81	24.43	53.81	66.96	8.92
iSEARLE-FFHQ	22.91	43.70	24.68	54.10	67.71	10.74
iSEARLE-VGGFace2	23.20	43.72	24.96	54.75	67.83	<u>10.97</u>
<b>iSEARLE-OTI</b>	<u>25.06</u>	<u>44.79</u>	<b>26.19</b>	<u>55.18</u>	<b>68.55</b>	10.94
<b>iSEARLE</b>	24.40	44.80	<u>25.23</u>	<b>55.69</b>	68.05	<b>11.24</b>

#### D. Additional Experiments

**Effect of  $\phi$  Pre-training Dataset** We carry out several experiments to study the impact of the  $\phi$  pre-training dataset. Specifically, besides the version of  $\phi$  trained on the test split of ImageNet1K, we also train some variants using: 1) CIRR training set, with 17K images; 2) FashionIQ training set, with 45K images; 3) NABirds [65] whole dataset, with 48K images depicting birds; 4) FFHQ [66] training set, with 50K images of aligned and cropped faces; 5) a subset of VGGFace2 [67] we obtained by randomly sampling 5 images per subject, resulting in 45K images of faces. Regarding CIRR and FashionIQ, we use only the raw images without considering the associated labels to keep the approach unsupervised. By relying on these two datasets, we assess the impact of pre-training  $\phi$  in a domain aligned with that of the testing dataset. Conversely, training on NABirds, FFHQ, or VGGFace2 provides insights into how effectively our method adapts to domains entirely distinct from the pre-training one. We selected these three datasets specifically for their very narrow domains, yet ensuring a sufficient number of images.

Table X reports the results. We notice how iSEARLE-FIQ achieves the best performance on FashionIQ and even outperforms iSEARLE-OTI, thus highlighting the effectiveness of our distillation-based approach. This result shows that pre-training on images belonging to the same domain as that of the testing ones leads to a performance gain. Moreover, iSEARLE-FIQ also manages to generalize to a broader domain obtaining promising results on both CIRR and CIRCO. Despite relying only on 17K training images, iSEARLE-CIRR achieves noteworthy results, suggesting that our approach is effective even in a low-data regime. Finally, we observe that iSEARLE-FFHQ, iSEARLE-NABirds, and iSEARLE-VGGFace2 obtain promising results despite being pre-trained on datasets that have highly specific domains that are extremely different from those of the testing datasets. To contextualize, Pic2Word [18] scores a  $R@1$  of 23.90 on CIRR and a  $mAP@10$  of 9.51 on CIRCO. In comparison, iSEARLE-VGGFace2, despite using a smaller backbone and being trained on such a narrow domain as human faces, achieves a  $R@1$  of 24.96 on CIRR and a  $mAP@10$  of 10.97 on CIRCO. These results show that our approach is robust to the  $\phi$  pre-training dataset. Moreover, our model demonstrates noteworthy generalization capabilities, making it well-suited for application in any real-world scenario

TABLE XI

EVALUATION OF THE VISUAL INFORMATION EMBEDDED IN  $v_*$  FOR DIFFERENT REGULARIZATION TECHNIQUES ON CIRR VALIDATION SET. IR AND CIR STAND FOR IMAGE RETRIEVAL AND COMPOSED IMAGE RETRIEVAL, RESPECTIVELY.

Ablation	Method	IR			CIR		
		R@1	R@3	R@5	R@1	R@5	R@10
OTI	w/o GPT reg	99.63	<b>100</b>	<b>100</b>	21.04	50.99	64.21
	random reg	99.26	99.95	99.95	<u>21.21</u>	48.43	62.97
	w/o reg	<b>99.77</b>	<b>100</b>	<b>100</b>	17.43	46.85	60.65
	<b>iSEARLE-OTI</b>	<b>99.77</b>	<b>100</b>	<b>100</b>	<b>25.57</b>	<b>57.11</b>	<b>70.46</b>
$\phi$	w/o reg	<b>99.21</b>	<b>99.95</b>	<u>99.95</u>	24.63	56.20	69.22
	w/o $\mathcal{L}_{pen}$	98.98	<b>99.95</b>	<b>100</b>	<u>25.66</u>	<u>57.14</u>	<u>70.22</u>
	<b>iSEARLE</b>	98.89	<b>99.95</b>	<u>99.95</u>	<b>25.74</b>	<b>57.35</b>	<b>70.32</b>

TABLE XII

COMPARISON WITH SUPERVISED BASELINES ON CIRR AND FASHIONIQ VALIDATION SETS. FOR FASHIONIQ, WE CONSIDER THE AVERAGE RECALL.

Method	FashionIQ		CIRR		
	R@10	R@50	R@1	R@5	R@10
Combiner-FIQ [4]	<b>32.96</b>	<b>54.55</b>	19.88	48.05	61.11
Combiner-CIRR [4]	20.91	40.40	<b>32.24</b>	<b>65.46</b>	<b>78.21</b>
<b>iSEARLE-OTI</b>	<b>25.06</b>	44.79	<b>25.57</b>	<b>57.11</b>	<b>70.46</b>
<b>iSEARLE</b>	24.40	<b>44.80</b>	<b>25.74</b>	<b>57.35</b>	70.32

without the requirement for domain-specific pre-training.

**Visual Information in  $v_*$**  We carry out an image retrieval experiment by studying whether the pseudo-word tokens can retrieve the corresponding images. In this way, we assess the effectiveness of the pseudo-word tokens in capturing visual information. Starting from an image  $I$ , we obtain the corresponding pseudo-word token  $v_*$  and its associated pseudo-word  $S_*$  through textual inversion. We craft a generic prompt including the pseudo-word  $S_*$ , such as “a photo of  $S_*$ ”. Then, we extract the text features via the CLIP text encoder and use them to query an image database. We expect the image  $I$  to be the top-ranked result if the pseudo-word token manages to effectively embed its visual content.

Tab. XI shows the results for Image Retrieval (IR) alongside the corresponding ones for Composed Image Retrieval (CIR). We use the CIRR validation set to conduct the experiments. We report the results of all the ablation studies related to the regularization technique for both OTI and  $\phi$ . Refer to Sec. V-C for more details on the setting of each ablation. Regardless of the regularization strategy, we observe that  $v_*$  effectively captures the visual information of the image, achieving almost perfect IR scores. However, we obtain a significant performance improvement in CIR when relying on our GPT-powered loss. This highlights how our regularization technique enhances the ability of the pseudo-word tokens to interact with the actual words composing the relative caption while preserving the visual information embedded in  $v_*$ .

**Comparison with Supervised Baselines** We measure the generalization capabilities of supervised CIR models by performing a comparison with our zero-shot approach. In particular, we consider Combiner [4], which fuses image and

text CLIP features through a combiner network. We chose Combiner as we believe it represents the most similar method to ours among the supervised ones, as we both rely on an out-of-the-box CLIP model. We train two versions of Combiner based on the B/32 backbone on the FashionIQ and CIRR training sets, respectively, using the official repository. We test both Combiner versions on FashionIQ and CIRR validation sets and report the results in Tab. XII. As expected, Combiner achieves the best performance when the training and testing datasets correspond. However, both supervised models struggle to generalize to different domains, as also observed by [18]. Conversely, iSEARLE exhibits remarkable performance on both datasets in a zero-shot manner. Thus, given that we do not require a costly manually annotated training set, the proposed method demonstrates better scalability and suitability for real-world applications of composed image retrieval.

## VI. CONCLUSION

In this work we expand upon our conference paper and introduce a new task, Zero-Shot Composed Image Retrieval (ZS-CIR), aimed at tackling CIR without requiring an expensive labeled training dataset. Since its introduction, several works have addressed ZS-CIR, highlighting its significance and relevance to the research community. We present an approach, named iSEARLE, that involves pre-training a lightweight textual inversion network via a distillation loss to retain the expressiveness of an optimization-based method while achieving a substantial efficiency gain. In addition, we introduce an open-domain benchmarking dataset for CIR, named CIRCO. CIRCO is the first CIR dataset featuring multiple labeled ground truths, reduced false negatives, and a semantic categorization of the queries. iSEARLE achieves state-of-the-art performance on FashionIQ, CIRR and the proposed CIRCO. Moreover, the proposed approach demonstrates better generalization capabilities than competing methods, as shown by two additional evaluation settings, namely object composition and domain conversion.

## ACKNOWLEDGMENTS

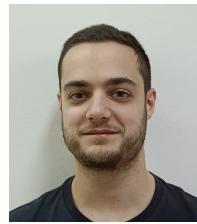
This work was partially supported by the European Commission under European Horizon 2020 Programme, grant number 101004545 - ReInHerit.

## REFERENCES

- [1] N. Vo, L. Jiang, C. Sun, K. Murphy, L.-J. Li, L. Fei-Fei, and J. Hays, “Composing text and image for image retrieval—an empirical odyssey,” in *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 6439–6448. 1, 2, 6
- [2] Z. Liu, C. Rodriguez-Opazo, D. Teney, and S. Gould, “Image retrieval on real-life images with pre-trained vision-and-language models,” in *Proc. of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 2125–2134. 1, 2, 6, 7, 8, 9
- [3] A. Baldrati, M. Bertini, T. Uricchio, and A. Del Bimbo, “Conditioned and composed image retrieval combining and partially fine-tuning clip-based features,” in *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 4959–4968. 1, 2, 9
- [4] ———, “Effective conditioned and composed image retrieval combining CLIP-based features,” in *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 21466–21474. 1, 2, 15

- [5] ——, “Composed image retrieval using contrastive learning and task-oriented clip-based features,” *ACM Transactions on Multimedia Computing, Communications and Applications*, vol. 20, no. 3, pp. 1–24, 2023. 1, 2, 3, 10
- [6] H. Wen, X. Zhang, X. Song, Y. Wei, and L. Nie, “Target-guided composed image retrieval,” in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 915–923. 1, 2
- [7] G. Delmas, R. S. Rezende, G. Csurka, and D. Larlus, “ARTEMIS: Attention-based retrieval with text-explicit matching and implicit similarity,” in *Proc. of International Conference on Learning Representations (ICLR)*, 2022. 1
- [8] S. Lee, D. Kim, and B. Han, “Cosmo: Content-style modulation for image retrieval with text feedback,” in *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 802–812. 1
- [9] A. Baldi, L. Agnolucci, M. Bertini, and A. Del Bimbo, “Zero-shot composed image retrieval with textual inversion,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 15 338–15 347. 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 13
- [10] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark et al., “Learning transferable visual models from natural language supervision,” in *Proc. of International Conference on Machine Learning (ICML)*. PMLR, 2021, pp. 8748–8763. 2, 3
- [11] R. Gal, Y. Alaluf, Y. Atzmon, O. Patashnik, A. H. Bermano, G. Chechik, and D. Cohen-Or, “An image is worth one word: Personalizing text-to-image generation using textual inversion,” in *Proc. of International Conference on Learning Representations (ICLR)*, 2023. 2, 3, 4
- [12] T. L. Berg, A. C. Berg, and J. Shih, “Automatic attribute discovery and characterization from noisy web data,” in *Proc. of the European Conference on Computer Vision (ECCV)*. Springer, 2010, pp. 663–676. 2, 6
- [13] X. Han, Z. Wu, P. X. Huang, X. Zhang, M. Zhu, Y. Li, Y. Zhao, and L. S. Davis, “Automatic spatially-aware fashion concept discovery,” in *Proc. of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2017, pp. 1463–1471. 2, 6
- [14] H. Wu, Y. Gao, X. Guo, Z. Al-Halah, S. Rennie, K. Grauman, and R. Feris, “Fashion iq: A new dataset towards retrieving images by natural language feedback,” in *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 11 307–11 317. 2, 6, 7, 8, 9
- [15] X. Guo, H. Wu, Y. Cheng, S. Rennie, G. Tesauro, and R. Feris, “Dialog-based interactive image retrieval,” in *Proc. of Advances in Neural Information Processing Systems (NeurIPS)*, vol. 31, 2018. 2, 6
- [16] M. Forbes, C. Kaeser-Chen, P. Sharma, and S. Belongie, “Neural naturalist: Generating fine-grained image comparisons,” in *Proc. of the Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019, pp. 708–717. 2, 6
- [17] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *Proc. of the European Conference on Computer Vision (ECCV)*. Springer, 2014, pp. 740–755. 2, 6, 10, 12
- [18] K. Saito, K. Sohn, X. Zhang, C.-L. Li, C.-Y. Lee, K. Saenko, and T. Pfister, “Pic2word: Mapping pictures to words for zero-shot composed image retrieval,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 19 305–19 314. 2, 3, 6, 10, 11, 12, 13, 14, 15
- [19] V. W. Liang, Y. Zhang, Y. Kwon, S. Yeung, and J. Y. Zou, “Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 17 612–17 625, 2022. 2, 4
- [20] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, “Vqa: Visual question answering,” in *Proc. of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2015, pp. 2425–2433. 2
- [21] Z. Shao, Z. Yu, M. Wang, and J. Yu, “Prompting large language models with answer heuristics for knowledge-based visual question answering,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 14 974–14 983. 2
- [22] M. Cornia, M. Stefanini, L. Baraldi, and R. Cucchiara, “Meshed-memory transformer for image captioning,” in *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 10 578–10 587. 2
- [23] M. Barraco, S. Sarto, M. Cornia, L. Baraldi, and R. Cucchiara, “With a little help from your own past: Prototypical memory networks for image captioning,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 3021–3031. 2
- [24] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 10 684–10 695. 2, 3
- [25] C. Meng, R. Rombach, R. Gao, D. Kingma, S. Ermon, J. Ho, and T. Salimans, “On distillation of guided diffusion models,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 14 297–14 306. 2, 3
- [26] M. Levy, R. Ben-Ari, N. Darshan, and D. Lischinski, “Data roaming and early fusion for composed image retrieval,” *arXiv preprint arXiv:2303.09429*, 2023. 3, 9
- [27] L. Ventura, A. Yang, C. Schmid, and G. Varol, “Covr: Learning composed video retrieval from web video captions,” *arXiv preprint arXiv:2308.14746*, 2023. 3
- [28] Y. Liu, J. Yao, Y. Zhang, Y. Wang, and W. Xie, “Zero-shot composed text-image retrieval,” *arXiv preprint arXiv:2306.07272*, 2023. 3
- [29] J. Chen and H. Lai, “Pretrain like you inference: Masked tuning improves zero-shot composed image retrieval,” *arXiv preprint arXiv:2311.07622*, 2023. 3
- [30] W. Li, H. Fan, Y. Wong, M. Kankanhalli, and Y. Yang, “CAT-LLM: Context-aware training enhanced large language models for multi-modal contextual image retrieval,” 2024. 3, 9
- [31] S. Karthik, K. Roth, M. Mancini, and Z. Akata, “Vision-by-language for training-free compositional image retrieval,” in *The Twelfth International Conference on Learning Representations*, 2024. 3, 9
- [32] S. Sun, F. Ye, and S. Gong, “Training-free zero-shot composed image retrieval with local concept reranking,” *arXiv preprint arXiv:2312.08924*, 2023. 3
- [33] Y. Tang, J. Yu, K. Gai, Z. Jiamin, G. Xiong, Y. Hu, and Q. Wu, “Context-i2w: Mapping images to context-dependent words for accurate zero-shot composed image retrieval,” *arXiv preprint arXiv:2309.16137*, 2023. 3, 6, 10, 11, 13
- [34] G. Gu, S. Chun, W. Kim, Y. Kang, and S. Yun, “Language-only efficient training of zero-shot composed image retrieval,” *arXiv preprint arXiv:2312.01998*, 2023. 3, 10, 11, 12, 13
- [35] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell et al., “Language models are few-shot learners,” in *Proc. of Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, 2020, pp. 1877–1901. 3, 5
- [36] N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein, and K. Aberman, “Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 22 500–22 510. 3
- [37] N. Kumari, B. Zhang, R. Zhang, E. Shechtman, and J.-Y. Zhu, “Multi-concept customization of text-to-image diffusion,” in *Proc. of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3
- [38] N. Cohen, R. Gal, E. A. Meirom, G. Chechik, and Y. Atzmon, ““This is my unicorn, Fluffy”: Personalizing frozen vision-language representations,” in *Proc. of the European Conference on Computer Vision (ECCV)*, 2022. 3, 4, 5, 6, 10, 11, 12, 13
- [39] B. Korbar and A. Zisserman, “Personalised clip or: how to find your vacation videos,” in *Proc. of British Machine Vision Association (BMVA)*, 2022. 3
- [40] P. Sharma, N. Ding, S. Goodman, and R. Soricut, “Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning,” in *Proc. of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2018, pp. 2556–2565. 3, 11
- [41] G. Hinton, O. Vinyals, and J. Dean, “Distilling the knowledge in a neural network,” in *NIPS Deep Learning and Representation Learning Workshop*, 2015. 3
- [42] M. Mistretta, A. Baldi, M. Bertini, and A. D. Bagdanov, “Improving zero-shot generalization of learned prompts via unsupervised knowledge distillation,” in *European Conference on Computer Vision*. Springer, 2024, pp. 459–477. 3
- [43] L. Beyer, X. Zhai, A. Royer, L. Markeeva, R. Anil, and A. Kolesnikov, “Knowledge distillation: A good teacher is patient and consistent,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10 925–10 934. 3
- [44] G. Chen, W. Choi, X. Yu, T. Han, and M. Chandraker, “Learning efficient object detection models with knowledge distillation,” in *Proc. of Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, 2017. 3

- [45] A. Chawla, H. Yin, P. Molchanov, and J. Alvarez, "Data-free knowledge distillation for object detection," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 3289–3298. 3
- [46] A. Sauer, D. Lorenz, A. Blattmann, and R. Rombach, "Adversarial diffusion distillation," *arXiv preprint arXiv:2311.17042*, 2023. 3
- [47] M. Mistretta, A. Baldorati, L. Agnolucci, M. Bertini, and A. D. Bagdanov, "Cross the Gap: Exposing the Intra-modal Misalignment in CLIP via Modality Inversion," in *The Thirteenth International Conference on Learning Representations*, 2025. 4
- [48] S. Gu, C. Clark, and A. Kembhavi, "I can't believe there's no images! learning visual tasks using only language supervision," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 2672–2683. 4, 5
- [49] O. Tov, Y. Alaluf, Y. Nitzan, O. Patashnik, and D. Cohen-Or, "Designing an encoder for stylegan image manipulation," *ACM Transactions on Graphics (TOG)*, vol. 40, no. 4, pp. 1–14, 2021. 5
- [50] J. Zhu, Y. Shen, D. Zhao, and B. Zhou, "In-domain gan inversion for real image editing," in *Proc. of the European Conference on Computer Vision (ECCV)*. Springer, 2020, pp. 592–608. 5
- [51] A. Kuznetsova, H. Rom, N. Alldrin, J. Uijlings, I. Krasin, J. Pont-Tuset, S. Kamali, S. Popov, M. Malluci, A. Kolesnikov *et al.*, "The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale," *International Journal of Computer Vision (IJCV)*, vol. 128, no. 7, pp. 1956–1981, 2020. 5, 10
- [52] D. Hendrycks and K. Gimpel, "Gaussian error linear units (GELUs)," *arXiv preprint arXiv:1606.08415*, 2016. 5
- [53] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *Proc. of International Conference on Machine Learning (ICML)*. PMLR, 2020, pp. 1597–1607. 5
- [54] J. D. Robinson, C.-Y. Chuang, S. Sra, and S. Jegelka, "Contrastive learning with hard negative samples," in *International Conference on Learning Representations*, 2020. 6
- [55] Y. Kalantidis, M. B. Sarıyıldız, N. Pion, P. Weinzaepfel, and D. Larlus, "Hard negative mixing for contrastive learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 21798–21809, 2020. 6
- [56] J. MacQueen *et al.*, "Some methods for classification and analysis of multivariate observations," in *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, vol. 1, no. 14. Oakland, CA, USA, 1967, pp. 281–297. 6
- [57] R. Gal, M. Arar, Y. Atzman, A. H. Bermano, G. Chechik, and D. Cohen-Or, "Encoder-based domain tuning for fast personalization of text-to-image models," *ACM Transactions on Graphics (TOG)*, vol. 42, no. 4, pp. 1–13, 2023. 6, 14
- [58] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International Journal of Computer Vision (IJCV)*, vol. 115, pp. 211–252, 2015. 6, 10, 12
- [59] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778. 9
- [60] D. Hendrycks, S. Basart, N. Mu, S. Kadavath, F. Wang, E. Dorundo, R. Desai, T. Zhu, S. Parajuli, M. Guo *et al.*, "The many faces of robustness: A critical analysis of out-of-distribution generalization," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 8340–8349. 10, 12
- [61] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," in *International Conference on Learning Representations*, 2018. 10
- [62] J. Yu, Z. Wang, V. Vasudevan, L. Yeung, M. Seyedhosseini, and Y. Wu, "Coca: Contrastive captioners are image-text foundation models," *Transactions on Machine Learning Research*, vol. Aug 2022, 2022. 10
- [63] Z. Wang, A. Chen, F. Hu, and X. Li, "Learn to understand negation in video retrieval," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 434–443. 11
- [64] M. Yuksekgonul, F. Bianchi, P. Kalluri, D. Jurafsky, and J. Zou, "When and why vision-language models behave like bags-of-words, and what to do about it?" in *The Eleventh International Conference on Learning Representations*, 2022. 11
- [65] G. Van Horn, S. Branson, R. Farrell, S. Haber, J. Barry, P. Ipeirotis, P. Perona, and S. Belongie, "Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 595–604. 14
- [66] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 4401–4410. 14
- [67] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman, "Vggface2: A dataset for recognising faces across pose and age," in *2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)*. IEEE, 2018, pp. 67–74. 14



**Lorenzo Agnolucci** received the M.S. degree (cum laude) in Computer Engineering from the University of Florence, Italy, in 2021. Currently, he is a PhD student at the University of Florence at the Media Integration and Communication Center (MICC). His research interests revolve around machine learning and computer vision, with a particular focus on low-level vision and vision-language models.



**Alberto Baldorati** received the M.S. degree (cum laude) in Computer Engineering from the University of Florence, Italy, in 2021. Currently, he is a Ph.D. student enrolled in the Italian National PhD Program in AI at the University of Pisa, while actively conducting research at the Media Integration and Communication Center (MICC) affiliated with the University of Florence. His research interests include machine learning and computer vision, focusing on vision and language, composed image retrieval, and fashion image generation.



**Alberto Del Bimbo** (Senior Member, IEEE) received the master's degree cum laude in electrical engineering, in 1977. He is a Full Professor of computer engineering, and the Director of the Media Integration and Communication Center, University of Florence, Florence, Italy. From 1996 to 2000, he was the President of the IAPR Italian Chapter and from 1998 to 2000, the Member-at-Large with the IEEE Publication Board. His research interests include multimedia information retrieval, pattern recognition, and computer vision. Prof. Del Bimbo received the SIGMM Technical Achievement Award for Outstanding Technical Contributions to Multimedia Computing, Communications and Applications. He was nominated the ACM Distinguished Scientist in 2016. He is a co-founder of Small Pixels, an academic spin-off working on visual quality improvement based on AI.



**Marco Bertini** is an associate professor of computer science at the School of Engineering of the University of Florence and director of the Media Integration and Communication Center (MICC) at the same university. His interests regard computer vision, multimedia, pattern recognition, and their application to different domains such as cultural heritage. He has been general co-chair, program co-chair and area chair of several international conferences and workshops on multimedia and computer vision (ACM MM, ICMR, CBMI, etc.), and was associate editor of IEEE Transactions on Multimedia. He has been involved in different roles in more than 10 EU research projects. He is a co-founder of Small Pixels, an academic spin-off working on visual quality improvement based on AI.