

SCOT: Self-Supervised Contrastive Pretraining For Zero-Shot Compositional Retrieval

Bhavin Jawade^{1,2*}, João V. B. Soares¹, Kapil Thadani¹, Deen Dayal Mohan¹,
Amir Erfan Eshratifar¹, Benjamin Culpepper¹, Paloma de Juan¹,
Srirangaraj Setlur², Venu Govindaraju²
¹Yahoo Research, ²University at Buffalo, SUNY

bhavinja@buffalo.edu, jvbsoares@yahooinc.com, thadani@yahooinc.com
deendayal.mohan@yahooinc.com, erfan.eshratifar@yahooinc.com, jackcul@yahooinc.com
pdjuan@yahooinc.com, setlur@buffalo.edu, govind@buffalo.edu

Abstract

Compositional image retrieval (CIR) is a multimodal learning task where a model combines a query image with a user-provided text modification to retrieve a target image. CIR finds applications in a variety of domains including product retrieval (e-commerce) and web search. Existing methods primarily focus on fully-supervised learning, wherein models are trained on datasets of labeled triplets such as FashionIQ and CIRR. This poses two significant challenges: (i) curating such triplet datasets is labor intensive; and (ii) models lack generalization to unseen objects and domains. In this work, we propose SCOT (Self-supervised COMpositional Training), a novel zero-shot compositional pretraining strategy that combines existing large image-text pair datasets with the generative capabilities of large language models to contrastively train an embedding composition network. Specifically, we show that the text embedding from a large-scale contrastively-pretrained vision-language model can be utilized as proxy target supervision during compositional pretraining, replacing the target image embedding. In zero-shot settings, this strategy surpasses SOTA zero-shot compositional retrieval methods as well as many fully-supervised methods on standard benchmarks such as FashionIQ and CIRR. Our code and models are available at <https://github.com/yahoo/SCOT>.

1. Introduction

The field of image retrieval is advancing rapidly, with growing interest in multimodal queries that incorporate both images and text. Compositional Image Retrieval (CIR) is a recently proposed task that aims at retrieving images us-

¹*Work done during research internship at Yahoo Research.

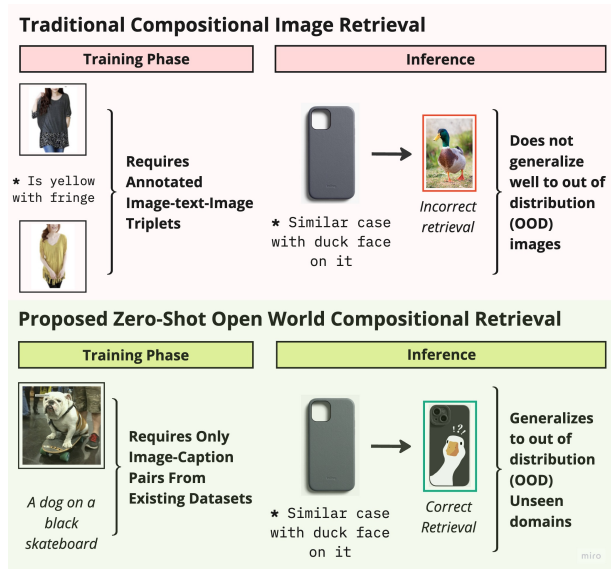


Figure 1. Compositional image retrieval methods typically require domain-specific image-text-image triplets for training and cannot generalize to unseen domains. In contrast, SCOT uses existing large noisy captioned image datasets for compositional training and demonstrates zero-shot generalizability to new domains.

ing a query composed of both an image and text [14, 37]. The query or reference image defines some initial desired elements, while the text describes the relative modification that a user would like to see in the retrieved images. CIR provides users with a versatile way to communicate their intent through iterative query refinement, which is potentially valuable in a broad range of real-world tasks such as product retrieval in e-commerce and fine-grained web search.

CIR can be framed as a multimodal fused representation learning task in which the goal is to train an effective feature fusion network. This sits in contrast to other well-studied

vision-language tasks such as image-text matching, image captioning, and visual question answering, as CIR uniquely learns a representation to jointly capture visual cues and text descriptors that match the target image of interest. Most CIR methods [9, 12, 19, 37] are trained in a fully-supervised manner using curated human-annotated datasets of triplets, with each triplet consisting of a reference image, a user-provided modification text, and a target image.

Current supervised CIR approaches do not generalize well to unseen domains or zero-shot scenarios, as illustrated in Fig 1. They are dependent on the availability of large datasets of image-text-image triplets, which are typically domain-specific and have limited applicability to open-world settings. Manual labeling for new triplet datasets is also labor-intensive. To overcome these challenges, a recent line of work explores zero-shot CIR using textual inversion [2, 6, 31, 35], e.g., using image-text pairs to learn to map images into text token embeddings. An image-derived token embedding—which can be thought of as corresponding to a *pseudo-token*—can then be combined with text token embeddings from the modification text and encoded as text to produce a composite embedding for retrieval. These approaches do not require annotated image-text-image triplets and can adapt to new domains thanks to the generalizability of contrastively-pretrained image-text encoders.

In this work, we propose a novel pretraining strategy for zero-shot CIR (ZS-CIR) which we name SCOT (Self-supervised COmpositional Training). This approach does not require human-annotated triplets and demonstrates open-world generalizability by using captioned images from large and varied datasets. We specifically exploit the proximity of visual and textual representations of the same concept in the embedding space of large-scale contrastively-pretrained vision-language models, which enables the use of target text embeddings instead of target image embeddings for supervision. Given an image and its caption, we first generate a training example by feeding the caption into a large language model (LLM) and prompting it to output a creative modification text and a corresponding modified caption. A CIR model is then trained by using the reference image and the generated modification text as input, with the generated modified caption as the target.

SCOT models are trained to compose reference images with modification texts by optimizing a contrastive image retrieval loss. This differs from inversion-based techniques, which do not directly train a composition model but rely on the composition capabilities of existing frozen pretrained image-text encoders. SCOT pretraining is agnostic to the choice of composition model, which can include unfrozen encoders [25] and early image-text fusion [21]. Comprehensive experiments show that SCOT surpasses current ZS-CIR techniques and nears fully-supervised performance on FashionIQ [40] and CIRR [26] without domain-specific

training. The key contributions of this work are:

1. We introduce a novel compositional pretraining strategy that requires only image-text pairs, using LLMs to create image-text-text triplets and pretrained vision-language models to encode both images and text.
2. We demonstrate zero-shot generalizability on domain-specific (FashionIQ [40]) and open-world (CIRR [26]) compositional retrieval datasets, showing that SCOT outperforms existing zero-shot approaches.
3. Through quantitative and qualitative experiments, we evaluate the impact of various parameters such as training dataset size, sample distribution, backbone and supervision type on zero-shot generalizability.

2. Related Work

Compositional Image Retrieval (CIR): Numerous methods have been proposed to learn composite representations of visual and text features for retrieval. Most research lies within the supervised setting [3, 4, 8, 9, 12, 17, 19, 21, 26, 27, 36, 37], with earlier work relying on fashion datasets containing human-annotated triplets [15, 40]. The DCNet approach [19] jointly trains feature extractors with a composition and correction network on FashionIQ [40]. CoSMo [20] uses content and style modulator networks to combine the image and text representations. FashionVLP [12] is a recently-proposed multimodal Transformer trained with a variety of fashion image inputs including crops, landmarks and ROIs. The need to go beyond fashion products and motivate research in open-world interactive retrieval led to the creation of open-domain annotated datasets: CIRR [26] (using images from NLVR2 [34]), CIRCO [2], and LaSCo [21] (the latter two using images from MS-COCO [24]). Despite this progress, the zero-shot generalizability of traditional fully-supervised models has been limited.

Zero-Shot Compositional Retrieval (ZS-CIR): To overcome these limitations, recent work [2, 5, 6, 13, 18, 25, 31, 35] has developed zero-shot annotation-free strategies for CIR. One line of work [2, 6, 31, 35] adopts textual inversion, which had previously found success in the text-to-image generation literature [11]. Recently, Saito et al. [31] proposed Pic2Word wherein an MLP is trained to map a picture to a pseudo-token, which the text encoder can then combine with the modification text to produce a composite embedding. Baldrati et al. [2] present SEARLE, which involves a two-stage process for training a textual inversion network. The first stage runs Optimization-based Textual Inversion (OTI) with CLIP [30] image and text encoders to find a text token embedding that corresponds to a given image encoding. In the second stage, those token embeddings are used as targets to learn a textual inversion network. Note that textual inversion approaches focus on learning how to invert the image into token embeddings, while taking advantage of the existing composition capabilities of pretrained

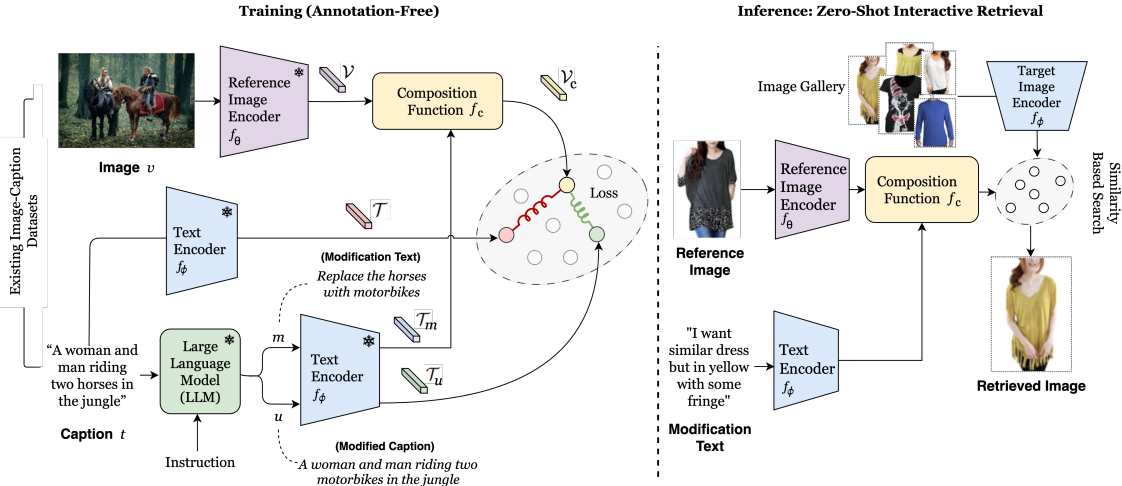


Figure 2. **SCOT pretraining and inference.** **Left:** The composition function f_c is trained using existing image-caption datasets, a frozen image-text encoder (such as CLIP), and a frozen large language model (LLM). The LLM generates the modification text m and a modified caption u . The reference image embedding \mathcal{V} and the modification text embedding \mathcal{T}_m are passed to f_c to get the composed embedding \mathcal{V}_c . We optimize the parameters of f_c to draw \mathcal{V}_c towards the modified caption \mathcal{T}_u and away from the original caption \mathcal{T} . The full loss also pushes \mathcal{V}_c away from the embeddings of other (non-matching) modified captions within each batch (not illustrated here). **Right:** During inference, we compute the similarity between the composed embedding and the embeddings of gallery images to retrieve the target image.

text encoders. In contrast, our approach directly optimizes a contrastive loss by training with triplets that closely mimic those of the CIR task. It can thus use any choice of composition model (including unfrozen encoders), and can be easily fine-tuned further with domain-specific data. A variety of other ZS-CIR approaches have been recently proposed. Gu et al. [13] train a denoising Transformer for image-text composition on 18M synthetic images along with 2B captioned images from LAION [32]. Karthik et al. [18] introduce CIReVL, a training-free approach that involves captioning the reference image, modifying the caption using an LLM and retrieving the target image using the modified caption. Chen and Lai [5] propose masking-augmented contrastive pretraining for visual and textual encoders to recover masked visual information through text prompts. Jang et al. [17] train a model to generate the modification text given a pair of images. The model is sued for generating synthetic training data, resulting in a semi-supervised approach. In concurrent work, Liu et al. [25] propose an approach for automatic construction of image-text-image training triplets. They source captioned images from the LAION-COCO dataset [33] and use either text templates or LLMs to generate modification texts and corresponding modified captions. Modified captions are then used to retrieve images to serve as supervision targets. The authors note that this approach of retrieving supervision target images from a corpus can be problematic due to the eventual absence of suitable images and/or retrieval errors [25]. In Section 4.4, we show example triplets illustrating these issues and present a controlled experiment demonstrating

that SCOT’s use of semantically-relevant text targets significantly outperforms the use of retrieved image targets.

3. Method

This section describes SCOT, a ZS-CIR technique requiring only captioned image datasets. The approach is outlined in Fig. 2. We review contrastively pretrained image-text encoders in Section 3.1. Sections 3.2, 3.3 and 3.4 detail our pretraining strategy, loss function, and inference.

3.1. Large-Scale Contrastive Pretraining

Following previous work, we use image and text representations from large-scale contrastively-pretrained models: CLIP [30], BLIP [23] and BLIP-2 [22]. CLIP (Contrastive Language-Image Pretraining) [30] aims to jointly learn visual and textual representations that are semantically aligned. For a given image-caption pair (v_i, t_i) , let $\mathcal{V}_i = f_\theta(v_i)$ denote the normalized image embedding from image encoder f_θ and $\mathcal{T}_i = f_\phi(t_i)$ denote the normalized text embedding from text encoder f_ϕ . CLIP contrastively enforces high similarity between positive pairs $(\mathcal{V}_i, \mathcal{T}_i)$ and low similarity between negative pairs $(\mathcal{V}_i, \mathcal{T}_j)$, $\forall i \neq j$. This is implemented via a symmetric cross-entropy loss over the similarity scores of image and text embeddings \mathcal{V}_i and \mathcal{T}_j . The image-to-text part of the loss is defined as:

$$\mathcal{L}_{i2t} = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{\langle \mathcal{V}_i, \mathcal{T}_i \rangle / \kappa}}{\sum_{j=1}^N e^{\langle \mathcal{V}_i, \mathcal{T}_j \rangle / \kappa}} \quad (1)$$





Reference Image	Original Caption	Modification Text	Modified Caption
	Two women sit and pose with stuffed animals .	Change the stuffed animals to guitars.	Two women sit and pose with guitars.
	Two men are working on a bicycle on the side of the road	Replace the bicycle with a skateboard for a more urban scene.	Two men are riding a skateboard on the side of the road.
	Alcuff Upholstered Power Recliner Sofa, 87" W, Blue	Switch the blue color for an earthy green for a more restful look.	Alcuff Upholstered Power Recliner Sofa, 87" W, Green
	Motorcycle 3D Printed Hard Back Case Mobile Cover	Replace the motorcycle design with a cityscape design for an urban feel.	Cityscape 3D Printed Hard Back Case Mobile Cover

Figure 3. **LLM-generated text triplet samples**, showing appropriate modifications over different image domains.

where $\langle \cdot, \cdot \rangle$ is the dot product, N the batch size and κ the temperature parameter.

BLIP [23] and BLIP-2 [22] are other pretraining approaches that demonstrate strong performance on benchmarks for image-text retrieval. BLIP-2 employs a lightweight trainable Querying Transformer (Q-Former) module whose image features come from a frozen pre-trained CLIP encoder. The initial training stage performs representation learning by jointly optimizing three objectives that include an image-text contrastive loss as in CLIP.

3.2. Self-Supervised Compositional Pretraining

Our approach is primarily motivated by the fact that contrastively-pretrained models are able to align related visual and textual representations in the embedding space. This enables us to use the aligned textual representation as a proxy for an image representation, thereby eliminating the need for a target image during training. For inference, we can search across gallery images by encoding them using the the contrastively-paired visual encoder.

The goal of this method is to train a composition operation f_c to combine the representations from a user-provided image and modification text. We rely on contrastively-paired image and text encoders, denoted respectively as f_θ and f_ϕ . Given a captioned image dataset $D = \{(v_i, t_i)\}_{i=1}^M$, we compute the image embeddings $\mathcal{V}^i = f_\theta(v_i)$ and corresponding caption embeddings $\mathcal{T}^i = f_\phi(t_i)$.

We prompt a large language model (LLM) to generate a modification text m_i given an original caption t_i . The modification text will be used as one of the inputs to the composition function f_c during training. To provide the supervision signal for the predicted composed representation, we use the same LLM to generate a modified caption u_i , which should be similar to t_i but with modification m_i applied. Fig. 3 contains samples of LLM-generated triplets.

Next, we compute the embeddings for m_i and u_i using f_ϕ .

$$m_i, u_i \leftarrow \text{LLM}(t_i) \quad (2)$$

$$\mathcal{T}_m^i, \mathcal{T}_u^i = f_\phi(m_i), f_\phi(u_i) \quad (3)$$

We pass the modification text representation \mathcal{T}_m^i and the image representation \mathcal{V}^i through the learnable composition function f_c to obtain the composed image representation $\mathcal{V}_c^i = f_c(\mathcal{V}^i, \mathcal{T}_m^i)$. We use the recently proposed Combiner network [3] to implement f_c . Briefly, it performs a learnable weighted fusion of the image and text embeddings. We encourage readers to refer to [3] or the supplementary material of this paper for details on the Combiner network.

3.3. Training Objective

We minimize a modified contrastive loss in order to pull the predicted composed embedding \mathcal{V}_c^i towards the generated target text embedding \mathcal{T}_u^i for an input sample (v_i, t_i) while pushing it away from target text embeddings $\mathcal{T}_u^j, \forall j \neq i$ from other examples within its batch. Let $S(x, y)$ denote the cosine similarity between vectors x and y , i.e., $S(x, y) = \frac{x \cdot y}{\|x\|_2 \|y\|_2}$. We define:

$$\mathcal{L}_{\text{pos}} = -\log \sum_{i=1}^N e^{S(\mathcal{V}_c^i, \mathcal{T}_u^i)} \quad (4)$$

$$\mathcal{L}_{\text{neg}} = \log \sum_{i,j}^N e^{S(\mathcal{V}_c^i, \mathcal{T}_u^j) \cdot (1 - \delta_{ij})} \quad (5)$$

where δ_{ij} is the Kronecker delta function, which is 1 when $i = j$ and 0 otherwise.

Previous work in traditional cross-modal retrieval [10, 39] has demonstrated the effectiveness of hard-negative mining. To improve the robustness of our embeddings, we follow [10, 38] and adopt a *margin-based* hard-negative mining strategy. Let λ be a fixed scalar margin, then we define:

$$S_\lambda(x, y) = S(x, y) \cdot \Theta(S(x, y) > \lambda) \quad (6)$$

where Θ is the Heaviside step function, which is 1 if the condition inside is true and 0 otherwise. This is used in Eq. 5 resulting in an updated negative loss.

$$\mathcal{L}'_{\text{neg}} = \log \sum_{i,j}^N e^{S_\lambda(\mathcal{V}_c^i, \mathcal{T}_u^j) \cdot (1 - \delta_{ij})} \quad (7)$$

For stronger supervision, we also include the original unmodified caption embeddings $\mathcal{T}^j \forall j \leq N$ as hard negatives for \mathcal{V}_c^i . This moves the composed representation away from the original caption and closer to the desired modified caption, ensuring it does not retain features from the original sample that are absent in the target caption. We use all the

original captions in a batch as negatives for that batch, resulting in the following combined loss for negatives.

$$\mathcal{L}''_{\text{neg}} = \mathcal{L}'_{\text{neg}} + \log \sum_{i,j}^N e^{S_\lambda(\mathcal{V}_e^i, \mathcal{T}^j)} \quad (8)$$

Using Eqs. (4, 8) we minimize the following final loss with respect to the parameters of the composition function f_c :

$$\mathcal{L} = \alpha_{\text{pos}} \cdot \mathcal{L}_{\text{pos}} + \alpha_{\text{neg}} \cdot \mathcal{L}''_{\text{neg}} \quad (9)$$

where α_{pos} and α_{neg} are positive and negative scaling factors respectively. In summary, the composition function is trained to apply the LLM-generated modification text to the reference image such that the resulting composed representation lies close to the embedding of the modified caption.

3.4. Inference

As shown in Fig. 2, all gallery images for retrieval are encoded with the image encoder f_θ . During inference, we combine the embeddings of the reference image and user-provided modification text using the learned composition function f_c , as in training. This composite representation is used to retrieve the most similar gallery images by computing cosine similarity with their image embeddings.

4. Experiments

We now turn to quantitative and qualitative evaluations of SCOT for ZS-CIR. Additional results are in the appendix.

4.1. Datasets

We train on three datasets of captioned images: MSCOCO [24] (189K pairs), Flickr30K [41] (45K pairs), and ABO [7] (58K pairs), totaling 290K image-text pairs. Following previous works [2, 25, 31], we assess zero-shot capabilities on FashionIQ [40] and CIRRR [26], two compositional retrieval datasets with annotated triplets. Here, FashionIQ assess zero-shot generalizability in the fashion domain and CIRRR on open-world retrieval setting.

4.2. Implementation Details

Encoders. Unless otherwise stated, we use BLIP-2¹ as a frozen² image and text encoder.

Textual triplet generation. To generate modification texts m and modified captions u , we use the instruction-tuned Falcon-7B LLM [1]. As directly prompting this model produces noisy and inconsistent generations on our task, we generate 4K text triplets from the better-performing GPT-4 model [29] and use them for LoRA fine-tuning [16] of 4-bit quantized Falcon-7B [1]. Finally, we generate a dataset

¹We use BLIP-2 with EVA-CLIP ViT-G/14 backbone from LAVIS.

²While encoders can also be finetuned with SCOT, we keep them frozen to compare fairly with prior work, most of which uses frozen encoders.

of over 290K text triplets using the finetuned Falcon-7B, which can be reused in subsequent training runs. SCOT is not reliant on any specific LLM, so newer or stronger models can also be used to refine and expand the triplet dataset.

Other training details. We train with AdamW [28], batch size 1024 and learning rate 1×10^{-4} . In the loss (Sec. 3.3), we set positive scaling factor $\alpha_{\text{pos}} = 10$ and negative scaling factor $\alpha_{\text{neg}} = 0.1$, and margin $\lambda = 0.2$. For the Combiner, we use the same hyperparameters as the original work [3]. Training and inference uses 2 NVIDIA A100 GPUs.

4.3. Comparison with state-of-the-art methods

Evaluation metrics. We present a quantitative comparison against the state-of-the-art on the FashionIQ [40] and CIRRR [26] datasets. The evaluation metric for FashionIQ is the average recall at rank K ($R@K$). Following prior work [2, 25, 31] we present $R@10$ and $R@50$ on the validation set. For CIRRR, we follow the authors’ proposed protocol to report $\text{Recall}@K$ at four different ranks, i.e., $K \in \{1, 5, 10, 50\}$, along with $\text{Recall}_{\text{subset}}@K$, which uses small subsets with fully labelled negatives for each query image [26]. We show results for existing zero-shot approaches and fully-supervised approaches.

Baselines. As reference, we present results of retrieving using just the image embedding (Image-Only), just the modification text embedding (Text-Only), or the sum of the two (Image+Text). For a fair comparison against prior zero-shot methods such as Pic2Word [31] and SEARLE [2], which rely on frozen backbones, we include results from TransAgg [25] with frozen backbones. Baldrati et al. [2] present two variants of their approach: SEARLE-OTI, which requires inference-time optimization, and SEARLE, which trains a textual inversion network to reproduce the OTI outputs in a single forward pass. Here, we use the reported results for SEARLE and its larger version SEARLE-XL. Finally, we note that existing methods use different backbones, amounts and types of data, fusion architectures, and pretraining strategies. For instance, Pic2Word [31] uses 3M images with a frozen CLIP L/14 backbone within a textual inversion-based approach, whereas TransAgg [25] uses 32K synthetic triplets with BLIP and a Transformer-based fusion method. We provide results segregated by backbone in Tables 1 and 2, and further analyze the importance of different contrastively-trained backbones in Section 4.4.

Results on FashionIQ. From Table 1, the best-performing SCOT model improves by 11.78% on $R@10$ and by 13.8% on $R@20$ over SEARLE-XL [2]. SCOT also demonstrates notable data efficiency: utilizing only 290K image-text pairs for training, in contrast to the 3M images used in Pic2Word’s training with Conceptual Captions, we achieve 13.75% improvement over Pic2Word on $R@10$. The zero-shot performance of SCOT exceeds many fully-supervised methods, such as DCNet [19], CLIP4Cir [3], and Fashion-

	Backbone	Method	Average		Dress		Shirt		Top/Tee	
			R@10	R@50	R@10	R@50	R@10	R@50	R@10	R@50
Superv.	Multi	MAAF [9]	24.3	48.8	23.8	48.6	21.3	44.2	27.9	53.6
	Multi	DCNet [19]	30.44	58.29	28.95	56.07	23.95	47.30	30.44	58.29
	Multi	FashionVLP [12]	34.27	62.51	32.42	60.29	31.89	58.44	38.51	68.79
	CLIP L/14	CLIP4CIR [3]	38.32	61.74	33.81	59.40	39.99	60.45	41.41	65.37
	BLIP	BLIP4CIR [27]	43.49	67.31	42.09	67.33	41.76	64.28	46.61	70.32
Zero-Shot	CLIP B/32	Image-Only	5.88	13.19	6.96	14.08	4.46	11.89	6.22	13.61
		Text-Only	18.41	36.28	14.92	33.81	19.77	34.69	20.55	40.33
		Image+Text	13.36	27.51	12.44	28.55	12.61	24.82	15.04	29.16
		PALAVRA [6]	19.76	37.25	17.25	35.94	21.49	37.05	20.55	38.76
		SEARLE [2]	22.89	42.53	18.54	39.51	24.44	41.61	18.54	39.51
		TransAgg [25]	<u>23.91</u>	<u>44.68</u>	<u>19.44</u>	<u>42.04</u>	<u>25.37</u>	<u>42.69</u>	<u>26.93</u>	<u>49.31</u>
	SCOT (Ours)	24.14	43.44	19.73	41.24	25.51	42.93	27.18	46.14	
	CLIP L/14	Image-Only	7.97	17.43	5.25	13.63	10.54	20.65	8.10	18.01
		Text-Only	19.01	35.26	15.22	33.01	19.82	33.31	21.87	39.46
		Image+Text	18.12	33.17	14.27	31.33	19.13	32.28	20.95	35.90
		Pic2Word [31]	24.7	43.7	20.0	40.2	26.2	43.6	27.9	47.4
		SEARLE-XL [2]	25.56	46.23	20.48	43.13	26.89	45.58	29.32	49.97
		TransAgg [25]	28.57	48.29	23.85	44.57	29.54	47.79	32.33	52.52
	SCOT (Ours)	<u>28.27</u>	<u>47.44</u>	<u>23.69</u>	<u>45.06</u>	<u>29.09</u>	<u>47.01</u>	<u>32.02</u>	<u>50.33</u>	
	BLIP	Image-Only	6.65	15.40	5.05	12.19	7.55	17.76	7.34	16.26
		Text-Only	24.01	42.73	20.03	39.96	24.63	41.02	27.38	47.22
		Image+Text	8.06	18.16	6.14	19.78	9.37	19.87	8.66	19.78
		TransAgg [25]	<u>26.95</u>	<u>46.10</u>	<u>21.67</u>	<u>41.89</u>	<u>28.07</u>	<u>45.63</u>	<u>31.11</u>	<u>50.79</u>
SCOT (Ours)	30.68	51.33	26.42	49.23	30.91	49.65	34.72	55.12		
BLIP-2	Image-Only	7.53	17.93	4.21	11.89	10.59	23.51	7.81	18.41	
	Text Only	24.68	43.59	20.77	41.64	25.95	42.83	27.33	46.31	
	Image+Text	<u>29.21</u>	<u>50.05</u>	<u>23.30</u>	<u>45.61</u>	<u>32.82</u>	<u>53.09</u>	<u>31.51</u>	<u>51.45</u>	
	SCOT (Ours)	38.45	60.03	32.78	55.91	41.42	61.09	41.15	63.10	

Table 1. **Results on FashionIQ.** Zero-shot results from our proposed approach compared against existing zero-shot methods (bottom) presented alongside some fully-supervised approaches (top). For fair comparisons, SEARLE results are from the inversion model and TransAgg results are using frozen backbones. See supplementary material for more results.

VLP [12], while approaching that of BLIP4CIR [27].

Results on CIRR. From Table 2, SCOT exhibits improvements of 12.58% at R@1 and 10.86% at R@5 over SEARLE-XL. We also see that Text-Only performance is significantly higher than Image-Only performance on CIRR, and that naively adding image features to text degrades performance. This is explained by a known shortcoming of CIRR —also noted in prior work [2]—that modification texts often describe the target image completely, with reference images providing no additional information.

4.4. Discussion

1. Qualitative analysis. In Fig. 4 (Top) we present zero-shot qualitative retrieval results on FashionIQ, illustrating domain-specific behavior. The figure shows that SCOT effectively composes images and text to retrieve the most accurate product image. The second row is particularly interesting: all methods retrieve a gray tank top, but only SCOT specifically retrieves one with the Adidas logo, which was also present in the reference image. We evaluate the qualitative performance on open-world images using CIRR in Fig. 4 (Bottom). As discussed earlier, often in CIRR the



Figure 4. **Qualitative retrieval results** on validation sets. **Top:** FashionIQ [40]. **Bottom:** CIRR [26]. A green box indicates the correctly retrieved image. For CIRR, the rightmost column illustrates the corresponding modality weight learned by SCOT for that example. (Best viewed in color.)

modification text can be informative enough to retrieve the correct target image. The learned dynamic scalar scores of the Combiner network are shown in the last column. In

	Method	Backbone	Recall@K				Recall _{subset} @K		
			K = 1	K = 5	K = 10	K = 50	K = 1	K = 2	K = 3
Superv.	Multi	MAAF [9]	10.31	33.03	48.30	80.06	21.05	41.81	61.60
	OSCAR	CIRPLANT [26]	19.55	52.55	68.39	92.38	39.20	63.03	79.49
	CLIP L/14	CLIP4CIR [3]	33.59	65.35	77.35	95.21	62.39	81.81	92.02
	BLIP	BLIP4CIR [27]	40.15	73.08	83.88	96.27	72.10	88.27	95.93
Zero-Shot	CLIP B/32	Image-only	6.94	22.94	33.71	59.18	21.06	41.01	60.34
		Text-only	21.16	45.35	57.40	81.06	62.26	81.08	90.75
		Image+Text	10.46	32.41	46.39	75.11	30.09	54.24	73.20
		PALAVRA [6]	16.62	43.49	58.51	83.95	41.61	65.30	80.94
		SEARLE [2]	24.00	53.42	66.82	89.78	54.89	76.60	88.19
		TransAgg [25]	24.46	53.61	67.54	89.81	<u>57.81</u>	<u>78.17</u>	<u>89.54</u>
	SCOT (Ours)	22.80	53.18	66.22	89.64	53.25	75.45	88.31	
	CLIP L/14	Image-only	7.47	23.88	34.07	57.57	20.87	41.95	61.13
		Text-only	22.00	45.79	57.57	79.59	61.71	80.26	90.43
		Image+Text	10.55	32.70	45.71	74.26	31.06	55.69	73.93
		Pic2Word [31]	23.9	51.7	65.3	87.8	-	-	-
		SEARLE-XL [2]	24.24	52.48	66.29	88.84	53.76	75.01	88.19
		TransAgg [25]	25.04	53.98	67.59	<u>88.94</u>	<u>55.33</u>	<u>76.82</u>	<u>88.94</u>
	SCOT (Ours)	<u>24.36</u>	<u>53.52</u>	<u>67.37</u>	89.35	51.47	74.24	87.90	
	BLIP	Image-only	7.23	25.78	37.35	62.34	20.60	40.96	61.35
		Text-only	34.19	61.68	71.74	87.83	72.34	87.97	94.79
Image+Text		8.24	28.96	41.23	68.07	23.64	45.35	66.29	
TransAgg [25]		34.89	64.75	76.24	92.22	<u>66.34</u>	<u>83.76</u>	<u>92.92</u>	
SCOT (Ours)	36.31	66.19	77.37	92.96	<u>64.73</u>	83.20	92.15		
BLIP-2	Image-only	7.59	24.43	35.56	61.42	20.74	40.67	61.08	
	Text-only	<u>33.52</u>	<u>61.50</u>	<u>71.35</u>	<u>88.31</u>	<u>72.53</u>	<u>88.02</u>	94.87	
	Image+Text	19.69	49.98	64.39	90.01	45.69	71.18	85.83	
	SCOT (Ours)	36.82	64.34	74.48	93.42	75.73	88.70	<u>94.84</u>	

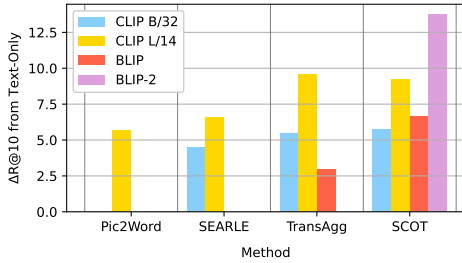
Table 2. **Results on CIRR.** Zero-shot results from our proposed approach compared against existing zero-shot methods (bottom) presented alongside some fully-supervised approaches (top). For fair comparisons, SEARLE results are from the inversion model and TransAgg results are using frozen backbones. See supplementary material for more results.

cases where the modification text completely describes the target image—such as in the third row—SCOT assigns a high weight to the text representation. In last row, it can be observed that the dog breed can only be inferred through the reference image; consequently SCOT assigns nearly equal weight to both image and text representations.

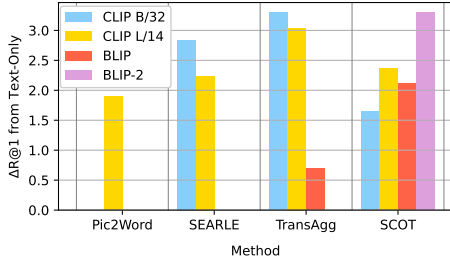
2. Impact of image-text alignment backbones. Using text embeddings as a proxy for image embeddings requires the image and text embedding spaces to be well-aligned. Here, we study the behavior of SCOT and other methods as we vary the encoder backbones. To recall, based on previous results [22, 23, 30], the relative ranking of the backbones we experimented with is CLIP-B/32 < CLIP-L/14 < BLIP < BLIP-2. From Table 1, with CLIP B/32, SCOT gets an average R@10 of 24.14% on FashionIQ [40]. With CLIP L/14, we observe 28.27%, nearly 4% higher. With BLIP, we observe another 2% improvement at R@10, while TransAgg produces a *drop* of 1.6%. Finally, for SCOT with BLIP-2, we see the largest improvement, of 8% over BLIP. On CIRR, as seen in Table 2, when using the CLIP B/32 backbone, SCOT is behind both TransAgg and SEARLE. SCOT then surpasses SEARLE when using the CLIP L/14 backbone, and surpasses TransAgg when switching to the BLIP

backbone. Thus, as with FashionIQ, the relative performance of methods changes with different backbones, with SCOT’s advantage increasing as the backbones improve. This can be more clearly seen in Fig. 5 where we present the relative gains of different methods with respect to the ‘Text-Only’ baseline. We define relative gain $\Delta R@K$ as the difference between Recall@K of a given method and that of the ‘Text-Only’ baseline with the corresponding backbone. On both FashionIQ and CIRR, Fig. 5 shows that as we improve the backbones, the relative gain of SCOT increases. Thus, not only does SCOT benefit from better backbones as represented by the performance of the ‘Text-Only’ baseline on those backbones, but its gain over that baseline also increases. Of note, with the BLIP backbone, SCOT has relative gains that are 2-3 larger than that of TransAgg with BLIP, showing that SCOT is unique in obtaining higher relative gains with better backbones.

3. Impact of dataset distribution. In Fig. 6, we illustrate how performance changes as we expand the training set. On both FashionIQ and CIRR, recall increases when utilizing larger subsets of the 189K MSCOCO image-caption pairs. This trend continues with the addition of Flickr30K. While both MSCOCO and Flickr30K contain generic real-world



(a) FashionIQ: $\Delta R@10$ from Text-Only



(b) CIRr: $\Delta R@1$ from Text-Only

Figure 5. **Gains relative to Text-Only.** Difference in recall (ΔR) between methods and the backbone-matched Text-Only baseline.

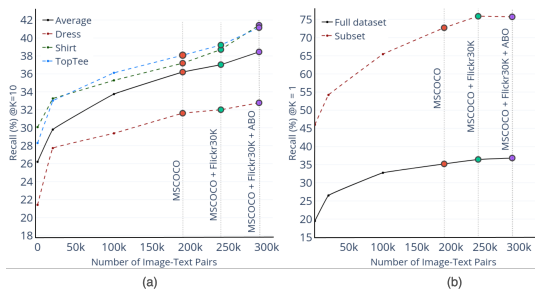


Figure 6. **Performance when changing the size and distribution of the training set,** evaluated on the (a) FashionIQ [40] validation set across clothing types and (b) CIRr [26] test set.

images, we wanted to also evaluate improvements brought by including domain-specific images. The Amazon Berkeley Objects (ABO) [7] dataset contains a variety of retail products, such as phone cases and furniture, accompanied by detailed captions. By including 58K image-caption pairs from ABO, we see around a 1% improvement on FashionIQ’s average $R@10$, going from 37.21% to 38.45%. Specifically for Shirt and Top/Tee, performance improves by around 2% when adding ABO. On CIRr, as shown in Fig. 6(b), incorporating ABO yields only a marginal gain in $R@1$ and no improvement in $R_{\text{subset}}@1$, likely due to the differing image distributions between ABO and CIRr.

4. Text supervision vs retrieved image supervision. An alternative way of using LLM-generated text triplets for ZS-CIR involves using each of the generated modified captions as a query to retrieve an image from a large corpus. Each retrieved image is then used as target supervision for

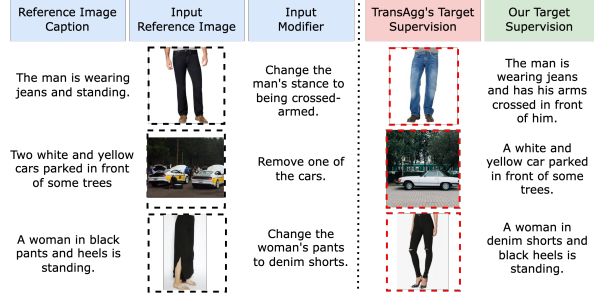


Figure 7. **Examples from LAION-CIR-LLM [25]** illustrating the challenges of using retrieved images as target supervision.

Supervision	Average		Dress		Shirt		Top/Tee	
	R@10	R@50	R@10	R@50	R@10	R@50	R@10	R@50
Image [25]	29.02	50.49	22.65	45.06	33.21	53.28	31.20	53.13
Text (Ours)	35.17	56.16	29.54	50.96	36.45	57.26	39.52	60.27

Table 3. **FashionIQ results with different supervision targets** when training on LAION-CIR-LLM [25] with a BLIP-2 backbone. We observe that the use of text targets for supervision performs significantly better than the image targets available in the dataset.

its corresponding reference image and generated modification text. Concurrently to our work, Liu et al. [25] experimented with this type of approach. Fig. 7 displays examples from the LAION-CIR-LLM dataset they proposed, which is based on image-caption pairs taken from LAION-COCO [33].³ As shown in the figure, the retrieved target images often do not match the expected modified caption due the absence of a relevant image in the corpus and/or retrieval errors. Table 3 presents an experiment comparing the use of the retrieved image target supervision from LAION-CIR-LLM [25] against text supervision using the dataset’s modified captions. The experiment uses BLIP-2 [22] as image and text encoder, and the Combiner [3] as composition function. We see that using retrieved images as targets gives an average $R@10$ on FashionIQ of 29.02%, whereas using text targets as proposed in our approach achieves 35.17%.

5. Conclusion

We propose a novel approach towards annotation-free ZS-CIR which leverages existing large captioned image datasets, along with contrastively-pretrained vision-language models. We demonstrate the zero-shot generalizability of this technique through extensive experimentation on domain-specific and open-world datasets. Our proposed approach, SCOT, achieves state-of-the-art performance in zero-shot settings while being on par with various fully-supervised approaches. We further substantiate this work with qualitative and quantitative experiments to analyze the impact of various components of our pretraining strategy.

³LAION-COCO (and by extension LAION-CIR-LLM) contains many clothing and product images, resulting in good coverage over FashionIQ.

References

- [1] Almazrouei, E., Alobeidli, H., Alshamsi, A., Cappelli, A., Cojocaru, R., Debbah, M., Goffinet, É., Hesselwood, D., Launay, J., Malartic, Q., et al.: The Falcon series of open language models. arXiv preprint arXiv:2311.16867 (2023) [5](#)
- [2] Baldrati, A., Agnolucci, L., Bertini, M., Del Bimbo, A.: Zero-shot composed image retrieval with textual inversion. arXiv preprint arXiv:2303.15247 (2023) [2](#), [5](#), [6](#), [7](#)
- [3] Baldrati, A., Bertini, M., Uricchio, T., Del Bimbo, A.: Effective conditioned and composed image retrieval combining CLIP-based features. In: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 21434–21442 (2022). <https://doi.org/10.1109/CVPR52688.2022.02080> [2](#), [4](#), [5](#), [6](#), [7](#), [8](#)
- [4] Barbany, O., Huang, M., Zhu, X., Dhua, A.: Leveraging large language models for multimodal search. In: CVPR Workshop on Fine-Grained Visual Categorization. pp. 1201–1210 (2024) [2](#)
- [5] Chen, J., Lai, H.: Pretrain like you inference: Masked tuning improves zero-shot composed image retrieval. arXiv preprint arXiv:2311.07622 (2023) [2](#), [3](#)
- [6] Cohen, N., Gal, R., Meirom, E.A., Chechik, G., Atzmon, Y.: “This is my unicorn, Fluffy”: Personalizing frozen vision-language representations. In: European Conference on Computer Vision. pp. 558–577. Springer (2022) [2](#), [6](#), [7](#)
- [7] Collins, J., Goel, S., Deng, K., Luthra, A., Xu, L., Gundogdu, E., Zhang, X., Yago Vicente, T.F., Dideriksen, T., Arora, H., Guillaumin, M., Malik, J.: ABO: Dataset and benchmarks for real-world 3D object understanding. CVPR (2022) [5](#), [8](#)
- [8] Delmas, G., Rezende, R.S., Csurka, G., Larlus, D.: ARTEMIS: Attention-based retrieval with text-explicit matching and implicit similarity. In: International Conference on Learning Representations (ICLR) (2022) [2](#)
- [9] Dodds, E., Culpepper, J., Herdade, S., Zhang, Y., Boakye, K.: Modality-agnostic attention fusion for visual search with text feedback. arXiv preprint arXiv:2007.00145 (2020) [2](#), [6](#), [7](#)
- [10] Faghri, F., Fleet, D.J., Kiros, J.R., Fidler, S.: VSE++: Improving visual-semantic embeddings with hard negatives. arXiv preprint arXiv:1707.05612 (2017) [4](#)
- [11] Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A.H., Chechik, G., Cohen-or, D.: An image is worth one word: Personalizing text-to-image generation using textual inversion. In: The Eleventh International Conference on Learning Representations (2023), <https://openreview.net/forum?id=NAQvF08TcyG> [2](#)
- [12] Goenka, S., Zheng, Z., Jaiswal, A., Chada, R., Wu, Y., Hedau, V., Natarajan, P.: FashionVLP: Vision language Transformer for fashion retrieval with feedback. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 14105–14115 (June 2022) [2](#), [6](#)
- [13] Gu, G., Chun, S., Kim, W., Jun, H., Kang, Y., Yun, S.: CompoDiff: Versatile composed image retrieval with latent diffusion. arXiv preprint arXiv:2303.11916 (2023) [2](#), [3](#)
- [14] Guo, X., Wu, H., Cheng, Y., Rennie, S., Tesauro, G., Feris, R.: Dialog-based interactive image retrieval. In: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (eds.) Advances in Neural Information Processing Systems. vol. 31 (2018), https://proceedings.neurips.cc/paper_files/paper/2018/file/a01a0380ca3c61428c26a231f0e49a09-Paper.pdf [1](#)
- [15] Han, X., Wu, Z., Huang, P.X., Zhang, X., Zhu, M., Li, Y., Zhao, Y., Davis, L.S.: Automatic spatially-aware fashion concept discovery. In: Proceedings of the IEEE international conference on computer vision (ICCV). pp. 1463–1471 (2017) [2](#)
- [16] Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: LoRA: Low-rank adaptation of large language models. In: International Conference on Learning Representations (2022), <https://openreview.net/forum?id=nZeVKeeFYf9> [5](#)
- [17] Jang, Y.K., Kim, D., Meng, Z., Huynh, D., Lim, S.N.: Visual delta generator with large multi-modal models for semi-supervised composed image retrieval. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16805–16814 (2024) [2](#), [3](#)
- [18] Karthik, S., Roth, K., Mancini, M., Akata, Z.: Vision-by-language for training-free compositional image retrieval. arXiv preprint arXiv:2310.09291 (2023) [2](#), [3](#)
- [19] Kim, J., Yu, Y., Kim, H., Kim, G.: Dual compositional learning in interactive image retrieval.

- Proceedings of the AAAI Conference on Artificial Intelligence **35**(2), 1771–1779 (May 2021). <https://doi.org/10.1609/aaai.v35i2.16271>, <https://ojs.aaai.org/index.php/AAAI/article/view/16271> **2, 5, 6**
- [20] Lee, S., Kim, D., Han, B.: CoSMo: Content-style modulation for image retrieval with text feedback. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 802–812 (June 2021) **2**
- [21] Levy, M., Ben-Ari, R., Darshan, N., Lischinski, D.: Data roaming and early fusion for composed image retrieval. arXiv preprint arXiv:2303.09429 (2023) **2**
- [22] Li, J., Li, D., Savarese, S., Hoi, S.: BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. arXiv preprint arXiv:2301.12597 (2023) **3, 4, 7, 8**
- [23] Li, J., Li, D., Xiong, C., Hoi, S.: Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In: International Conference on Machine Learning. pp. 12888–12900. PMLR (2022) **3, 4, 7**
- [24] Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: Common objects in context. In: Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13. pp. 740–755. Springer (2014) **2, 5**
- [25] Liu, Y., Yao, J., Zhang, Y., Wang, Y., Xie, W.: Zero-shot composed text-image retrieval. In: 34rd British Machine Vision Conference (BMVC) 2023 (2023) **2, 3, 5, 6, 7, 8**
- [26] Liu, Z., Rodriguez-Opazo, C., Teney, D., Gould, S.: Image retrieval on real-life images with pre-trained vision-and-language models. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 2125–2134 (2021) **2, 5, 6, 7, 8**
- [27] Liu, Z., Sun, W., Hong, Y., Teney, D., Gould, S.: Bi-directional training for composed image retrieval via text prompt learning. arXiv preprint arXiv:2303.16604 (2023) **2, 6, 7**
- [28] Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: International Conference on Learning Representations (2019), <https://openreview.net/forum?id=Bkg6RiCqY75>
- [29] OpenAI: GPT-4 technical report. arXiv preprint arXiv:2303.08774 **abs/2303.08774** (2023), <https://arxiv.org/abs/2303.08774> **5**
- [30] Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021) **2, 3, 7**
- [31] Saito, K., Sohn, K., Zhang, X., Li, C.L., Lee, C.Y., Saenko, K., Pfister, T.: Pic2Word: Mapping pictures to words for zero-shot composed image retrieval. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 19305–19314 (June 2023) **2, 5, 6, 7**
- [32] Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., et al.: LAION-5B: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems* **35**, 25278–25294 (2022) **3**
- [33] Schuhmann, C., Köpf, A., Vencu, R., Coombes, T., Beaumont, R.: LAION COCO: 600M synthetic captions from LAION2B-en. <https://laion.ai/blog/laion-coco/> (2022), accessed: Nov. 2023 **3, 8**
- [34] Suhr, A., Zhou, S., Zhang, A., Zhang, I., Bai, H., Artzi, Y.: A corpus for reasoning about natural language grounded in photographs. In: Korhonen, A., Traum, D., Márquez, L. (eds.) Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics. pp. 6418–6428. Association for Computational Linguistics, Florence, Italy (Jul 2019). <https://doi.org/10.18653/v1/P19-1644>, <https://aclanthology.org/P19-16442>
- [35] Tang, Y., Yu, J., Gai, K., Jiamin, Z., Xiong, G., Hu, Y., Wu, Q.: Context-I2W: Mapping images to context-dependent words for accurate zero-shot composed image retrieval. arXiv preprint arXiv:2309.16137 (2023) **2**
- [36] Tian, Y., Newsam, S., Boakye, K.: Fashion image retrieval with text feedback by additive attention compositional learning. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 1011–1021 (2023) **2**
- [37] Vo, N., Jiang, L., Sun, C., Murphy, K., Li, L.J., Fei-Fei, L., Hays, J.: Composing text and image for image retrieval - an empirical odyssey. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 6439–6448 (2019) **1, 2**

- [38] Wang, J., Song, Y., Leung, T., Rosenberg, C., Wang, J., Philbin, J., Chen, B., Wu, Y.: Learning fine-grained image similarity with deep ranking. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1386–1393 (2014) [4](#)
- [39] Wei, J., Xu, X., Yang, Y., Ji, Y., Wang, Z., Shen, H.T.: Universal weighting metric learning for cross-modal matching. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 13005–13014 (2020) [4](#)
- [40] Wu, H., Gao, Y., Guo, X., Al-Halah, Z., Rennie, S., Grauman, K., Feris, R.: Fashion IQ: A new dataset towards retrieving images by natural language feedback. In: Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition. pp. 11307–11317 (2021) [2](#), [5](#), [6](#), [7](#), [8](#)
- [41] Young, P., Lai, A., Hodosh, M., Hockenmaier, J.: From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics* **2**, 67–78 (2014), <https://aclanthology.org/Q14-1006> [5](#)