

# Benchmarking Robustness of Text-Image Composed Retrieval

Anonymous authors

Paper under double-blind review

## Abstract

Text-image composed retrieval aims to retrieve the target image through the composed query, which is specified in the form of an image plus some text that describes desired modifications to the input image. It has recently attracted attention due to its ability to leverage both information-rich images and concise language to precisely express the requirements for target images. However, the robustness of these approaches against real-world corruptions or further text understanding has never been studied. In this paper, we perform the first robustness study and establish three new diversified benchmarks for systematic analysis of text-image composed retrieval against natural corruptions in both vision and text and further probe textural understanding. For natural corruption analysis, we introduce two new large-scale benchmark datasets, CIRRR-C and FashionIQ-C for testing in open domain and fashion domain respectively, both of which apply 15 visual corruptions and 7 textural corruptions. For textural understanding analysis, we introduce a new diagnostic dataset CIRRR-D by expanding the original raw data with synthetic data, which contains modified text to better probe textural understanding ability including numerical variation, attribute variation, object removal, background variation, and fine-grained evaluation.

## 1 Introduction

Text-image composed retrieval, known as composed image retrieval or text-guided image retrieval, attempts to retrieve an image of interest from gallery images through a composed query of a reference image and its corresponding modified text. As a single word (‘dog’) can correspond to thousands of images (dogs in different breeds, poses, and scenarios), language is considered to be discrete and sparse, while images are regarded as dense and continuous. Using both images and text as queries enables the effective utilization of the continuous and dense nature of images to accurately express the requirements while leveraging discrete and sparse text to bridge semantic gaps beyond what the images alone can capture. It becomes fascinating for its potential in a wide range of real-world applications including fashion domain e-commerce Han et al. (2022b; 2023); Goenka et al. (2022); Han et al. (2022a); Chen et al. (2020) and open domain internet search Liu et al. (2021); Baldrati et al. (2022); Gu et al. (2023); Saito et al. (2023). However, existing text-image composed retrieval methods mostly are trained and tested on clean data, while the models in the real world may naturally encounter distribution shift Wang et al. (2021c), such as text typos and image corruptions owing to weather change. Additionally, there is currently no analysis of whether the model understands the meaning of the text rather than solely relying on finding correspondences with the main objects as a shortcut to the text-image composed retrieval task. For example, with a source image of a dog and modified text ‘change to two dogs on the table’, the model may retrieve the target image by only recognizing the word ‘dog’ and ‘table’ without the ability of numerical counting. Whether text-image composed retrieval models are robust in real-world applications, where natural corruption exists in both images and text, remains unexplored. Also, whether these models are robust across diverse textural understanding requirements, remains a domain blank.

In this work, we make the first attempt to evaluate the robustness of text-image composed retrieval by building three new large-scale robustness benchmarks on both fashion and open domains. We raise the following two questions: **Q1:** *How robust are text-image composed retrieval models to natural corruption including both visual and textual?* Further to evaluate the text understanding ability, we have the second question **Q2:** *How robust are text-image composed retrieval models to text understanding?*

To answer the first question, we present two benchmark datasets on text-image composed retrieval task. Based on two widely used datasets FashionIQ Wu et al. (2021) in the fashion domain and CIRRR Liu et al. (2021) in the open domain, we propose our benchmark datasets, namely FashionIQ-C and CIRRR-C, both with 15 visual corruptions and 7 textual corruptions to evaluate model robustness against natural corruption in both image and text. To answer the second question, we introduce a new diagnostic dataset CIRRR-D to probe text understanding abilities on five elementary scenarios including numerical variation, attribute manipulation, object removal, background variation, and fine-grained variation. In detail, we construct the diagnostic dataset by synthetic triplets based on the CIRRR validation set and use existing triplets from the CIRRR validation set with both main captions and extended captions. Our experiments show that the new benchmarks we introduced are suitable for robustness analysis against natural corruption on both image and text and further probe text understanding ability.

Our contributions are: (1) We make the first attempt to analyze the robustness of text-image composed retrieval methods against natural corruption (including visual and textual) and textual understanding (including five elementary variations). (2) We introduce three new large-scale benchmarks including two benchmark datasets (FashionIQ-C and CIRRR-C) to evaluate robustness against natural corruption in both image and text and one diagnostic benchmark CIRRR-D to probe text understanding robustness. (3) We present an empirical analysis and conduct extended experiments.

## 2 Related Works

**Robustness analysis.** Quantifying robustness aims to evaluate the model stability to defend against corruption including natural corruption Hendrycks & Dietterich (2019); Chantry et al. (2022); Wang et al. (2021b), adversarial attacks Croce et al. (2020); Wang et al. (2021a;b), or to probe certain ability such as logical reasoning Sanyal et al. (2022) and visual content manipulation Li et al. (2020). Traditional works about robustness analysis mainly focus on single modality involving visual modality-based tasks like image classification Hendrycks & Dietterich (2019), face detection Dooley et al. (2022); textual modality based task like text classification Zeng et al. (2021) and audio modality based task like speech recognition Mitra et al. (2017). Recently, robustness analysis against multimodal tasks, which is closer to real life and attempts to take a step towards a reliable system, has appeared but is still in its infancy. For example, Li et al. Li et al. (2020) take the first step to systematically analyze the robustness of a multimodal task, Visual Question Answering (VQA), against 4 generic robustness including linguistic variation and visual content manipulation. However, it is limited to VQA tasks and doesn't introduce benchmarks to pinpoint sophisticated reasoning abilities. Schiappa et al. Chantry et al. (2022) introduce natural corrupted visual and textual benchmarks on text-to-video retrieval. However, the robustness analysis of the multimodal underlying hypothesis, which aims to generalize textual semantic and reasoning ability to visual space, is not discussed. We consider the analysis of both natural corruption in image and text and further underlying text understanding and take the first step to conduct an extensive analysis of the natural corruption and text understanding of the robustness of deep neural networks in text-image composed retrieval.

**Diagnostic analysis.** Recently, a range of benchmarks for visual understanding have been proposed, including datasets for image captioning Shekhar et al. (2017), visual question answering Johnson et al. (2017), visual reasoning Zerroug et al. (2022) and visio-linguistic compositional reasoning Thrush et al. (2022); Yuksekgonul et al. (2022); Ma et al. (2023). For text-image composed retrieval, the benchmarks can be categorized into sythetic-based datasets by cubes Vo et al. (2019) or natural scenes Gu et al. (2023), fashion-based datasets Han et al. (2017); Berg et al. (2010); Wu et al. (2021), object-state dataset Isola et al. (2015) and open domain dataset Liu et al. (2021). Among them, the majority of the textual descriptions are limited by predefined attributes Han et al. (2017); Vo et al. (2019); Isola et al. (2015). To overcome this limitation, FashionIQ Wu et al. (2021) and CIRRR Liu et al. (2021) leverage the flexibility of natural language and becomes the most widely used benchmarks in fashion domain and open domain respectively. We expand and categorize the test set of CIRRR benchmark from its main and unused extended annotation to probe specific text understanding in numerical variation, attribute variation, object removal, background variation, and fine-grained variation. Similar to ours, many diagnostic datasets are also synthetic ones. CLVER Johnson et al. (2017) is a synthetic dataset to probe elementary vision reasoning including color, shape, and spatial relationships. CVR Zerroug et al. (2022) generates irregular shape, location, color, etc, and designed for



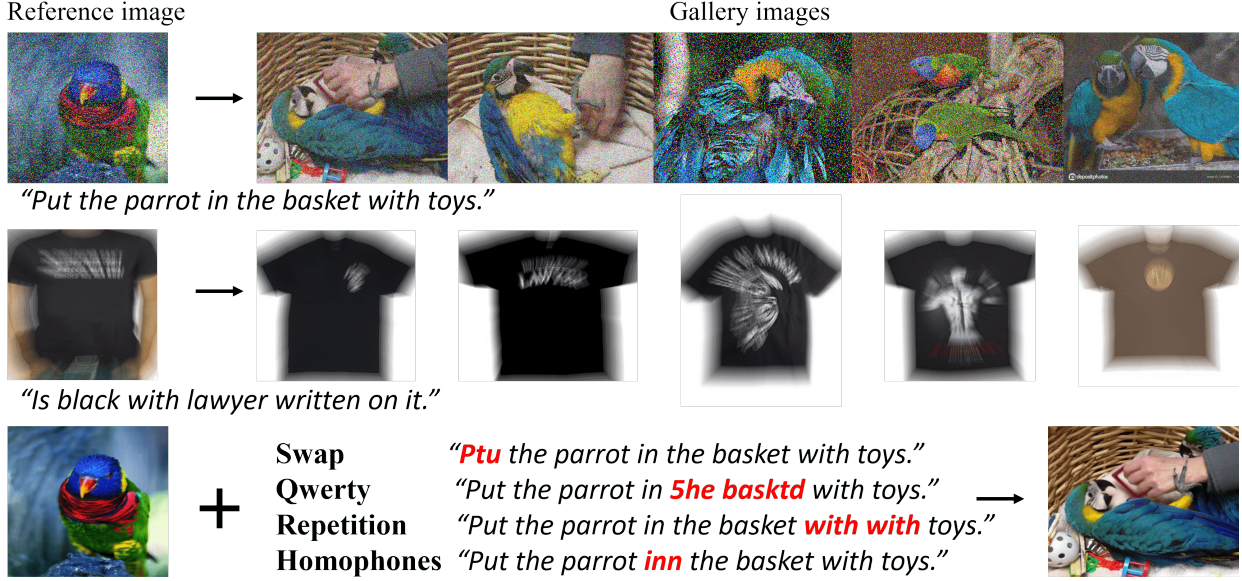


Figure 1: Sample visualization of proposed benchmarks under natural corruption with both visual and textual. **Top:** CIRRR-C with impluse noise image corruption **Middle:** FashionIQ-C with zoom blur image corruption; **Bottom:** CIRRR-C under character-level (Swap and Qwerty) and word-level (Repetition and Homophones) textual corruptions. Gallery images are shown without particular order.

detecting the outlier from a small set of generated images. However, these are all simulated images and are not generated by imitating natural scenes. CasualVQA Agarwal et al. (2020) and CompoDiff Gu et al. (2023) generate images imitating natural scenes. However, CasualVQA is designed for visual question answering tasks and the generated images include noticeable artifacts. While ComoDiff is designed for text-image composed retrieval, but they generate by replacing the objects (noun) only instead of attributes (numeral, adjective, manipulation instruction) like ours, so cannot pinpoint the target reasoning abilities. To detect the compositional abilities, visio-linguistic compositional reasoning diagnostic benchmarks supply variation of objects Ma et al. (2023); Yuksekgonul et al. (2022), attributes Ma et al. (2023); Yuksekgonul et al. (2022), or text order Yuksekgonul et al. (2022); Thrush et al. (2022). However, these benchmarks only supply image-text pairs for single-modality queries instead of image-text-image triplets for multi-modality queries. Also their text composition Ma et al. (2023); Yuksekgonul et al. (2022) may not have corresponding image like ‘grass eat horse’. In comparison, our CIRRR-D supplies the first visio-linguistic composition reasoning benchmark for text-image composed retrivel task in natural scenes.

**Text-image composed retrieval.** Composed image retrieval aims to retrieve the target image, where the input query is specified in the form of an image plus other interactions, such as relative attributeParikh & Grauman (2011), natural languageChen et al. (2020); Vo et al. (2019), spatial layoutMai et al. (2017), to describe the desired modifications. Among them, natural language as the most pervasive interaction between humans and computers to convey intricate specification has attracted increasing attention, which has often led ‘composed image retrieval’ to become interchangeable with ‘text-guided image retrieval’ in the literature. We term the task as text-image composed retrieval to clarify the composition of the query. Traditional text-image composed retrieval models implement separate independent image and text encoders, whose features are combined with late fusion. For example, TIRG Vo et al. (2019) and Artemis Delmas et al. (2022) implement separate pre-trained ResNet as image encoder and LSTM as text encoder. Until recently, with the power of unified multimodal space CLIP Radford et al. (2021), current text-image composed retrieval models achieved a noticeable improvement. For example, CLIP4CIR Baldrati et al. (2022) implements a light adapter as image-text late fusion and further tune it in target domains. Further based on CLIP, FAME Han et al. (2023) and CASE Levy et al. (2023) separately implement early cross attention between text and image, which shows obvious improvement.



Figure 2: Sample visualization of proposed benchmarks probing text understanding. **Top:** CIRRR-D with attribute variations; **Middle:** CIRRR-D with number variations; **Bottom:** CIRRR-D with object removal. CIRRR-D with background variation and fine-grained variation can be found in the supplementary materials.

### 3 Robustness Criteria for Text-Image Composed Retrieval

**Foundation of text-image composed retrieval.** Text-image composed retrieval aims to retrieve the target visual content through dense continuous images guided by sparse discrete text. We discuss ‘dense’ or ‘sparse’ in the semantic space, where a single semantic word can correspond to thousands of images. Therefore, text-image composed query can overcome the limitation of singular modality image retrieval, where text-image retrieval suffers from the unprecise description and unlimited correct targets, and image-image retrieval suffers from expression limitation without the ability to generalise to different visual content. In light of this, the foundation abilities of text-image composed retrieval are threefold: (1) Image representation to supply a precise anchor in the dense continuous visual space; (2) Text representation to supply subtle or significant differences between various visual contents, providing an unprecise target direction the model can generalize to; (3) Generalize sparse modified text attributes to dense reference images to precisely predict the target visual content by different variations of modality fusion.

**Definition of robustness in text-image composed retrieval.** According to the foundation of text-image composed retrieval above, a robust model should perform stable image feature extraction, text feature extraction, and modality fusion. In light of this, the robustness of text-image composed retrieval can be

defined in twofold: *robustness against natural corruption* for both text and image and *robustness against textual understanding* for consistent reasoning between textual and visual modality. Specifically, for robustness against natural corruption, we evaluate text-image composed retrieval models under ubiquitous corruption frequently encountered in real-life in both visual and textual. We evaluate 15 standard image corruptions with 5 severity categorized into noise, blur, weather, and digital following Hendrycks & Dietterich (2019). We also evaluate 7 text corruptions categorized into character-level and word-level. Further, for robustness against textual understanding, we evaluate common linguistic reasoning by selecting modified text with specific keywords or gallery set, categorized into numerical variation, attributes variation, object removal, background variation, and fine-grained variation respectively.

**Evaluation metrics.** To evaluate the performance of models in text-image composed retrieval, we adopt the standard evaluation metric in retrieval, i.e. Recall@K denoted as R@K for short. Further to measure robustness, we adopt relative robustness metrics  $\gamma = 1 - (R_c - R_p) / R_c$  following Chantry et al. (2022); Hendrycks & Dietterich (2019), where  $R_c$  and  $R_p$  are the R@K under clean data and data with corruption respectively. Additionally, in order to facilitate fair comparison among different models, we have expanded Delmas et al. (2022) and established a unified testing platform for the convenient integration of various models. In detail, we set gallery the whole validation set as in Delmas et al. (2022); Baldrati et al. (2022), which includes more distractors and results to higher discriminative requirement, instead of setting gallery the same as the query set as in Chen et al. (2020); Lee et al. (2021). Specifically for evaluating the fashionIQ dataset, we combine the two captions in a single query as Baldrati et al. (2022); Dodds et al. (2020) instead of combining the two modified captions in forward and reverse direction as Lee et al. (2021). All the evaluated models are trained in three categories jointly and tested individually for dress, shirt, and toptee categories. The reported results for fashionIQ are average of the three categories.

**Evaluation datasets.** We utilize three new benchmarks for our text-image composed retrieval experiments, which are generated from two existing datasets: FashionIQ Wu et al. (2021) in the fashion domain and CIRRR Liu et al. (2021) in the open domain. Both datasets include human-generated captions that distinguish image pairs. FashionIQ is based on the fashion domain containing 77,684 garment images, which can be divided into three categories: dress, shirt, and toptee. Each image in FashionIQ contains a single subject positioned centrally with a clean background. CIRRR is composed of 21,552 real-life images extracted from NLVR2 Suhr et al. (2018), which contains rich visual content in diverse backgrounds. As shown in Figure. 1, we build our benchmark and evaluate text-image composed retrieval models on text and image natural corruption robustness. Further as shown in Figure. 2, we expand the CIRRR dataset and evaluate text understanding robustness.

*Natural vision and text corruption.* To evaluate the robustness of the text-image composed retrieval model against natural corruption in both image and text, we create our robustness benchmark CIRRR-C and FashionIQ-C with 15 visual corruptions and 7 textual corruptions. For vision corruption, we follow Hendrycks & Dietterich (2019) to implement 15 standard natural corruptions which fall into four categories: noise, blur, weather, and digital, each having a severity from 1 to 5. For text corruption, we follow Rychalska et al. (2019) and implement the most related seven corruptions including four character-level corruptions and three word-level corruptions respectively.

*Diagnostic datasets.* Following the current methods Liu et al. (2021); Baldrati et al. (2022) reporting the results on the validation set, we expand and build our probing datasets CIRRR-D based on the validation set of CIRRR to pinpoint text understanding ability including numerical variation, attributes composition, object removal, background transformation, and fine-grained evaluation. We hypothesize that the model’s corresponding capabilities can be evaluated when the modified text involves descriptions such as numbers, attributes, objects removal, or changing the background; and the ability to deal with fine-grained variations can be evaluated when the gallery images are highly similar following Liu et al. (2021). In light of this, we build the triplets (reference image, modified text, and target image) of our probing dataset according to the appearances of specific keywords in modified text: "zero" to "ten", "number" for numerical query; color, shape and size for attribute query, "remove" for object removal query; "background" for background variation query. As shown in Table 4, the construction of CIRRR-D dataset can fall into five probing categories from three sources as follows: (1) The existing validation set of CIRRR is composed of 2297 images and 4181 triplets, which is currently widely used. Each image has a subset which is composed of 6 highly similar images as a

gallery to better detect fine-grained discriminative ability. (2) The auxiliary captions of the CIRR validation set, which is supplied but has not been used in the conventional evaluation. These captions indicate the differences in removal content or background changes between image pairs, but they may not provide enough information to accurately locate the target image. Therefore, we manually eliminated triplets that led to an excessive number of target images. (3) Synthetic images we generate through Visual ChatGPT Wu et al. (2023) to bridge the language reasoning and visual recognition. Based on the image from CIRR validation set, we generate its image caption by Visual ChatGPT and generate ten variants of the captions by ChatGPT including four for number variants, three for color variants, two for size variants, one for removing the objects respectively. Afterward, based on the reference image and caption variants, Visual ChatGPT utilizes groundingDINO Liu et al. (2023) to do object detection, segment anything Kirillov et al. (2023) to generate mask and stable diffusion Rombach et al. (2021) to generate target image. We manually eliminated the unplauble generated images. Our synthetic image pairs preserve the original background and only modify the specific areas mentioned in the text.

**Evaluated models.** We perform our experiments on six text-image composed retrieval models. The modality fusion of these approaches is roughly summarized in Supplementary Table 5 including Vo et al. (2019); Liu et al. (2021); Baldrati et al. (2022); Delmas et al. (2022); Han et al. (2022b); Dodds et al. (2020), which can be divided into overlapped categories: (1) Large pretrained model: FashionViL, CIRPLANT, CLIP4CIR, whose pretrained dataset size are 1.35 million, 6.5 million, and 400 million image-text pairs respectively; (2) Multi-task model: FashionViL, which is pretrained with four tasks at the same time; (3) Light attention-based methods: ARTEMIS; (4) Transformer-based models: MAAF, CIRPLANT, CIRPLANT and FashionViL. (5) Lightweight models: MAAF, TIRG, and ARTEMIS (all with ResNet50 image encoder and LSTM text encoder for fair comparison). (6) Single-modality models: Image-only (RN50) and Image-only (CLIP) are queried with images embedded by ResNet50 (same as evaluated TIRG, MAAF, ARTEMIS) and CLIP image encoder RN50x4 (same as evaluated CLIP4CIR) respectively. Text-only model is queried with text embedded using the CLIP text encoder. These methods were chosen because the reproduced results match with original reported results. We test FashionViL Han et al. (2022b) in the fashion domain, CIRPLANT in the open domain, and all the rest published models in both the fashion domain and open domain.

**Evaluation settings.** In order to ensure the fairness of the evaluation, we establish a standardized testbed for variant models except FashionViL to unify the evaluation process. Further to reproduce the original performance, we implement the official pre-trained weight for FashionViL and CLIP4CIR. We retrain the model and receive similar results as reported for models MAAF, TIRG, ARTEMIS, and CIRPLANT. In detail, we extend the existing ARTEMIS code framework to facilitate the convenient interface of different trained models, where TIRG and ARTEMIS are already implemented. For image input among these models, CIRPLANT is based on frozen ResNet152 pre-trained features while other models take raw images as input. We implement a frozen ResNet152 image encoder so that we can introduce corruption to the raw image directly. For all the 15 image corruptions, we perform the highest severity of corruption for obvious performance. For text input among these models, TIRG, MAAF, and ARTEMIS build the vocabulary based on appearance words in target evaluation caption, while FashionViL, CIRPLANT, and CLIP4CIR implement their vocabulary from large pretraining dataset. We implement textual corruption and change the raw text directly. For a concise explanation, we show the result of FashionIQ by showing the average of the three categories: dress, shirt, and toptee.

**Implementation details** We supply 5 models (TIRG, MAAF, ARTEMIS, CIRPLANT and CLIP4CIR) in the same testbed to have a fair comparison with different benchmark datasets. The selection of models or datasets can be easily accomplished through input parameters. Further models can be implemented in our testbed by simply providing model structure files with the necessary interface. In detail, the necessary interfaces include image feature extraction, text feature extraction, feature composing process and distance comparison. Our testbed is compatible with two environments and five models. CIRPLANT is implemented with Python(3.1) and Pytorch(1.8.1). TIRG, MAAF, CLIP4CIR and ARTEMIS were implemented with Python(3.8) and Pytorch(2.0). All the experiments are conducted and trained on NVIDIA A100 GPUs. We will maintain our code for benchmarking and testbed open source.



Table 1: Relative robustness score for text-image composed retrieval under 15 natural image corruptions in CIRR-C Recall@10 and FashionIQ-C Recall@10. Recall@10 performance under clean conditions on the left. **Bold** is the highest relative robustness for the five composed retrieval methods.

	Noise				Blur				Weather				Digital			
<b>CIRR-C</b>	Clean	Gauss.	Shot	Impluse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Bright	Contrast	Elastic	Pixel	JPEG
Image-only(RN50)	50.4	0.57	0.55	0.58	0.68	0.28	0.82	0.45	0.38	0.34	0.64	0.86	0.20	0.48	0.76	0.88
Image-only(CLIP)	36.2	0.56	0.55	0.58	0.66	0.32	0.83	0.49	0.52	0.45	0.77	0.91	0.24	0.41	0.78	0.91
Text-only(CLIP)	51.2	0.79	0.76	0.81	0.85	0.29	1.0	0.55	0.65	0.70	0.89	1.0	0.19	0.40	0.96	1.0
TIRG Vo et al. (2019)	55.1	0.34	0.36	0.34	0.48	0.21	0.70	0.43	0.31	0.22	0.40	0.70	0.12	0.47	0.74	0.84
MAAF Dodds et al. (2020)	49.9	0.50	0.49	0.50	0.62	0.26	0.80	0.41	0.36	0.31	0.50	0.74	0.11	0.48	0.83	0.87
ARTEMIS Delmas et al. (2022)	59.0	0.39	0.42	0.38	0.51	0.25	0.70	0.44	0.31	0.26	0.45	0.71	0.10	0.47	0.75	0.86
CIRPLANT Liu et al. (2021)	68.8	<b>0.70</b>	<b>0.69</b>	<b>0.71</b>	0.77	0.28	0.89	0.51	0.44	0.43	0.66	0.88	<b>0.17</b>	<b>0.56</b>	0.85	0.92
CLIP4CIR Baldrati et al. (2022)	<b>80.3</b>	0.68	0.68	0.69	<b>0.77</b>	<b>0.28</b>	<b>0.90</b>	<b>0.52</b>	<b>0.550</b>	<b>0.600</b>	<b>0.80</b>	<b>0.91</b>	0.16	0.39	<b>0.91</b>	<b>0.92</b>
<b>FashionIQ-C</b>	Clean	Gauss.	Shot	Impluse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Bright	Contrast	Elastic	Pixel	JPEG
TIRG Vo et al. (2019)	23.8	0.28	0.26	0.23	0.34	0.22	0.61	0.57	0.32	0.27	0.37	0.61	0.12	0.64	0.85	0.85
MAAF Dodds et al. (2020)	23.4	0.31	0.27	0.25	0.44	0.21	0.67	0.53	0.29	0.24	0.31	0.54	0.13	0.54	0.83	0.83
ARTEMIS Delmas et al. (2022)	24.9	0.24	0.24	0.20	0.38	0.26	0.65	0.60	0.36	0.25	0.38	0.55	0.14	0.63	0.86	0.87
FashionViL Han et al. (2022b)	23.4	0.26	0.28	0.25	0.40	<b>0.31</b>	<b>0.82</b>	<b>0.67</b>	0.33	0.31	0.34	<b>0.70</b>	0.15	<b>0.86</b>	<b>1.09</b>	<b>1.06</b>
CLIP4CIR Baldrati et al. (2022)	<b>35.9</b>	<b>0.44</b>	<b>0.42</b>	<b>0.44</b>	<b>0.54</b>	0.21	0.72	0.50	<b>0.460</b>	<b>0.430</b>	<b>0.60</b>	<b>0.70</b>	<b>0.22</b>	0.37	0.74	0.83

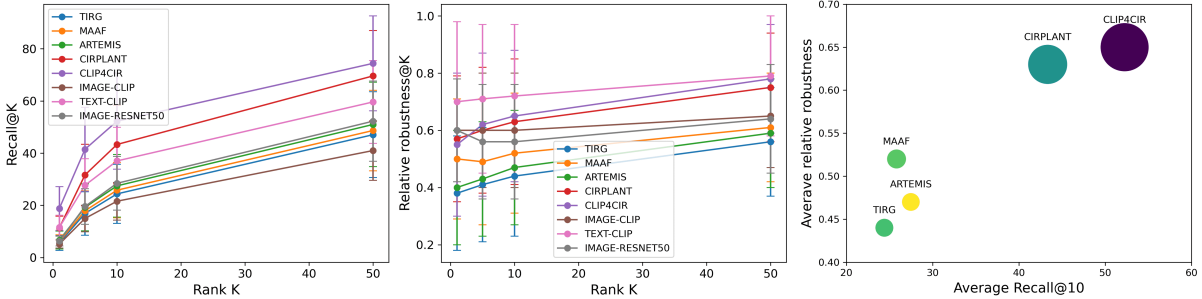


Figure 3: Models average performance in CIRR under 15 vision corruptions. Left: Recall vs. rank K. Middle:  $\gamma$  vs. rank K. Right:  $\gamma$  vs. recall@10, circle size indicates number of model parameters.

## 4 Results and Analysis

### 4.1 Natural corruption analysis

In order to evaluate whether the text-image composed retrieval models are robust under natural corruptions, we perform experiments under 15 visual corruptions which can be further categorized into **noise**, **blur**, **weather**, **digital** corruptions on both fashion domain and open domain. We report the relative robustness  $\gamma$  under the highest severity of each natural visual corruption in Table 1, which shares the same trend as other corruption severities. To evaluate the robustness against textual corruption, we conduct experiments under 7 textual corruptions on text-image composed retrieval models in **character-level** and **word-level**. Analysis about textual corruption is shown in supplementary 4.2.

**Pretraining.** Among the compared models, FashionViL CIRPLANT and CLIP4CIR are pretrained on large datasets, with respective sizes of 1.35 million, 6.5 million, and 400 million image-text pairs, while another three compared models are based on ImageNet pretrained ResNet50 as image encoder and random initialized LSTM as text encoder. As shown in Figure 3, the models with large pretrained datasets consistently show better robustness in both open domain and fashion domain. *This implies that models with large pretrained*

Table 2: Relative robustness score for text-image composed retrieval under 7 natural text corruptions in CIRR-C recall@10 and FashionIQ-C recall@10 on average of three categories. Recall@10 performance under clean conditions on the left. **Bold** are the highest relative robustness among the five compared methods.

	Character					Word		
<b>CIRR-C</b>	Clean	Swap	QWERTY	RemoveChar	RemoveSpace	Misspelling	Repetition	Homophone
Text-only	51.2	0.75	0.74	0.78	1.0	0.99	0.98	0.92
TIRG Vo et al. (2019)	55.1	0.77	0.76	0.80	1.0	0.98	1.0	0.89
MAAF Dodds et al. (2020)	49.9	<b>0.95</b>	<b>0.97</b>	<b>0.96</b>	<b>1.0</b>	<b>1.0</b>	<b>1.0</b>	<b>0.97</b>
ARTEMIS Delmas et al. (2022)	59.0	0.61	0.58	0.65	<b>1.0</b>	0.98	0.98	0.82
CIRPLANT Liu et al. (2021)	68.8	0.92	0.93	0.93	<b>1.0</b>	<b>1.0</b>	<b>1.0</b>	<b>0.97</b>
CLIP4CIR Baldrati et al. (2022)	<b>80.3</b>	0.89	0.89	0.90	<b>1.0</b>	<b>1.0</b>	0.99	<b>0.97</b>
	Character					Word		
<b>FashionIQ-C</b>	Clean	Swap	QWERTY	RemoveChar	RemoveSpace	Misspelling	Repetition	Homophone
TIRG Vo et al. (2019)	23.8	0.26	0.20	0.29	0.66	0.63	0.61	0.52
MAAF Dodds et al. (2020)	23.4	0.40	0.39	0.39	0.70	0.68	0.68	0.62
ARTEMIS Delmas et al. (2022)	24.9	0.25	0.20	0.31	0.70	0.67	0.67	0.55
FashionViL Han et al. (2022b)	23.4	<b>0.55</b>	<b>0.59</b>	<b>0.60</b>	<b>0.86</b>	<b>0.84</b>	<b>0.85</b>	<b>0.76</b>
CLIP4CIR Baldrati et al. (2022)	<b>35.9</b>	0.52	0.51	0.54	0.71	0.70	0.69	0.67

datasets may result in better robustness against visual corruptions, which is in alignment with the statement from Paul et al. Paul & Chen (2022).

**Bottleneck of robustness.** We visualize the recall performance with rank  $K$  improvement on the left of Figure 3 and relative robustness on the right. With ranking  $K$  improved, both the recall and relative robustness improved for all of the models. Additionally, we find out in Figure 3 that ResNet50-based image-only search shows even better accuracy as well as robustness compared with TIRG, ARTEMIS, and MAAF, which utilizes the same image encoder ResNet50. According to the three foundations of text-image composed retrieval in Sec. 3, both the image encoder and modality fusion module can be vulnerable to vision corruption. We observe that the ResNet50 backbone shows relatively high robustness while the modality fusion of the three models (TIRG, MAAF, and ARTEMIS) exacerbates the instability of the model. However, this phenomenon is not applicable to the CLIP feature, whose text and image embedding are aligned in a unified space in the pretraining process. Comparing Image-only (CLIP) and CLIP4CIR, which query with CLIP image embedding and CLIP text-image composed embedding respectively, we can find out CLIP4CIR consistently performs better recall performance as well as robustness in Figure 3. *Thus, we speculate that text features from aligned space can help boost the robustness, while text features from independent space will damage the model robustness.*

Further to better pinpoint the vulnerability in variant model fusion modules, we compare TIRG, MAAF and ARTEMIS with the same text LSTM and image backbone ResNet50 but different fusion methods. As shown in Table 1 in open domain, MAAF shows the trend performing the most robust, while ARTEMIS performs the second best over TIRG. The modality fusion modules of these three models, TIRG, ARTEMIS, and MAAF are concatenation-based, light attention, and transformer, respectively. Among them, MAAF utilizes modality-agnostic attention which extracts word and image tokens to conduct thorough merge through self-attention and cross-attention. *We hypothesize that more sufficient cross-modal interactions, such as cross-attention, can better promote robustness.*

## 4.2 Textual robustness against natural corruption

Comparing the relative robustness against textual corruption in Table 2 and visual corruption in Table 1, we can observe that the robustness is higher against textual corruptions. Among the compared models, MAAF and FashionViL show the highest robustness in open domain and fashion domain respectively. This aligns with the findings we discovered regarding vision corruptions, where large pretrained models (FashionViL

Table 3: Recall of CIRRR-D dataset. The red and green arrows indicate the performance increase or decrease compared with CIRRR queries. **Bold** and underline are the largest decrease and increase.

	R@5						Rsub@1
	CIRRR	Numerical	Attribute	Removal	Background	Fine grained	
Image-only(RN50)	31.55	31.47 ↓ (0.08)	32.57 ↑ (1.02)	35.99 ↑ (4.44)	39.15 ↑ (7.60)		20.25
Image-only(CLIP)	22.51	24.80 ↑ (2.29)	29.09 ↑ (6.58)	27.90 ↑ (5.39)	25.64 ↑ (3.13)		20.02
Text-only	39.02	42.84 ↑ (3.82)	49.45 ↑ (10.43)	11.62 ↓ (27.4)	11.62 ↓ (27.4)		53.73
TIRG Vo et al. (2019)	36.35	39.64 ↑ (3.29)	37.77 ↑ (1.42)	30.41 ↓ (5.94)	32.82 ↓ (3.53)		35.90
MAAF Dodds et al. (2020)	32.19	32.53 ↑ (0.34)	35.57 ↑ (3.38)	31.09 ↓ (1.10)	34.27 ↑ (2.08)		28.63
ARTEMIS Vo et al. (2019)	40.05	39.56 ↓ (0.49)	42.68 ↑ (2.63)	33.26 ↓ (6.79)	35.56 ↓ (4.49)		40.80
CIRPLANT Liu et al. (2021)	48.82	45.07 ↓ (3.75)	47.73 ↓ (1.09)	41.12 ↓ (7.70)	45.98 ↓ (2.84)		38.19
CLIP4CIR Baldrati et al. (2022)	62.94	64.18 ↑ (1.24)	69.15 ↑ (6.21)	31.66 ↓ (31.28)	41.88 ↓ (21.06)		62.66

with fashion-specific pretraining) and models with sufficient modality fusion result to higher robustness. Additionally, comparing CLIP4CIR and the Text-only retrieval method implementing CLIP text embedding, we find out the robustness can be boosted after fusion with vision modality. However, with images corrupted, CLIP4CIR shows lower robustness than text-only model, from which *we speculate that aligned clean image feature can boost the robustness, while the corrupted image feature will impair robustness.*

### 4.3 Robustness against text understanding

In this section, we analyze model reasoning ability through variation of modified text on numerical variation, attributes variation, object removal, background variation, and fine-grained variation, which are supplied by our proposed CIRRR-D dataset. We take the 4181 queries from the origin CIRRR dataset and evaluate them in our expanded CIRRR-D dataset as a baseline, which incorporates diverse reasoning instruction and can represent the models’ average performance across various instructions. However, unlike other function detection, the gallery set for fine-grained is a subset, which is composed of six highly similar images following Liu et al. (2021). We first evaluate performance on each query type as shown in Table 3. Detailed analyses are discussed below.

**Numerical variation.** To probe the ability of numerical variation, the modified text contains either a precise value of the number from zero to ten or an estimated value by comparison like ‘*reduce/increase the number*’. Comparing numerical specific query with CIRRR query as shown in Table 3, we can not observe significant variation which may result from the long-tailed distribution. Namely, the numerical set has a large number of samples in the range of 1 to 3, while a small number of samples in the range of 4 to 10. More analysis can be found in the supplementary D. For now, we speculate that *numerical modification may not be the bottleneck of the current text-image composed retrieval.*

**Attributes variation.** To evaluate the model’s discriminative ability when querying elementary attributes, the modified text includes variations of color, shape and size. As observed in Table 3, all of the methods (except CIRPLANT) achieve higher performance with attribute queries than with CIRRR queries. Additionally, the performance of CLIP based image-only model and CLIP4CIR have an obvious increment of over 6% compared with their performance with CIRRR queries, which have a strong ability of attribute recognition including color, shape, and size. This implies that *attribute is the one of main focuses during training and models gain strong attribute discriminative ability.*

**Object removal.** Object removal is a convenient approach to describe the differences between images but is universally overlooked by current methods in text-image composed retrieval. To probe the ability of object removal through CIRRR-D, the modified text of the query explicitly contains the word ‘*remove*’. As shown in Table 3, all of the five compared methods achieve their lowest performance in object removal with an average decrement of 10.6% compared with the CIRRR query. In particular, CLIP4CIR has a drop of over 30%, which may be a result of its static pretraining process by aligning only image text pairs without comparison between images. Surprisingly, image-only methods can have an increment over CIRRR queries, which illustrates that

visual similarity can boost the robustness over object removal but the text condition over guidance the model decision. This aligns with the foundation of the task: images are dense and continuous while text is sparse and discrete. *In the case of object removal, text guidance expands the possibility of the targets, which distracts the model and results in lower performance.*

**Background variation.** To probe the robustness against background modifications, the modified text of the query explicitly includes the word "background". We observe a similar phenomenon as in object removal, where the performance of compared models (except MAAF) decreases but the performance of image-only models increases compared to CIRRR queries. As the CIRRR-D sample visualization in Figure 6, we can observe that the background modification method is limited such as changing the background color or making the background blur, which can lead to unrelated targets by relying solely on the text itself. We further speculate that *a modified text leading to more satisfactory candidates may result in impaired outcomes.*

**Fine-grained variation.** To probe the fine-grained variation discriminative ability, we utilize the subset in the CIRRR dataset, where each image is retrieved from its subset composed of another five highly similar images. As the gallery is different from the above reasoning function, the recall cannot be compared with CIRRR query performance directly. We can observe from Table 3 that image-only models perform similarly to random guessing and text-only by CLIP embedding can achieve an acceptable result. Among the five compared methods, TIRG achieve the lowest performance which hypothesizes the slight adjustment in visual space is sufficient rather than exploring text deeply and establishing a fusion space. This phenomenon indicates that it is difficult to distinguish between two states in continuous visual space. In contrast, text can precisely define subtle differences due to its discrete nature. We also speculate that *a modified text offers accurate information while minimizing the number of feasible targets can enhance the model’s discriminative ability.*

## 5 Conclusion

In this work, we proposed three robustness benchmarks for text-image composed retrieval including two for natural corruption in both image and text and one for probing textual understanding. Concretely, we first introduced two benchmark datasets, CIRRR-C and FashionIQ-C with natural corruption (both image and text) in the open domain and fashion domain respectively. Further, we create benchmark CIRRR-D to assess the text understanding including number, attribute, object removal, background, and fine-grained variation. Based on our observation, we provide the following suggestions to enhance model robustness in text-image composed retrieval: 1) model pretrained on large datasets with little distribution shift will lead to better robustness, 2) text features from aligned space can help boost the robustness, while text features from independent space will damage the model robustness, 3) a modified text is more likely to enhance the model’s discriminative ability when it minimizes the number of feasible targets and will distract the model when it leads to more satisfactory candidates. These findings can potentially boost the robustness of text-image composed retrieval in the future.

## References

- Vedika Agarwal, Rakshith Shetty, and Mario Fritz. Towards causal vqa: Revealing and reducing spurious correlations by invariant and covariant semantic editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9690–9698, 2020.
- Alberto Baldrati, Marco Bertini, Tiberio Uricchio, and Alberto Del Bimbo. Effective conditioned and composed image retrieval combining clip-based features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21466–21474, 2022.
- Tamara L Berg, Alexander C Berg, and Jonathan Shih. Automatic attribute discovery and characterization from noisy web data. In *Computer Vision–ECCV 2010: 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5–11, 2010, Proceedings, Part I 11*, pp. 663–676. Springer, 2010.



- Madeline Chantry, Shruti Vyas, Hamid Palangi, Yogesh Rawat, and Vibhav Vineet. Robustness analysis of video-language models against visual and language perturbations. *Advances in Neural Information Processing Systems*, 35:34405–34420, 2022.
- Yanbei Chen, Shaogang Gong, and Loris Bazzani. Image search with text feedback by visiolinguistic attention learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3001–3011, 2020.
- Francesco Croce, Maksym Andriushchenko, Vikash Sehwal, Edoardo Debenedetti, Nicolas Flammarion, Mung Chiang, Prateek Mittal, and Matthias Hein. Robustbench: a standardized adversarial robustness benchmark. *arXiv preprint arXiv:2010.09670*, 2020.
- Ginger Delmas, Rafael Sampaio de Rezende, Gabriela Csurka, and Diane Larlus. Artemis: Attention-based retrieval with text-explicit matching and implicit similarity. *arXiv preprint arXiv:2203.08101*, 2022.
- Eric Dodds, Jack Culpepper, Simao Herdade, Yang Zhang, and Kofi Boakye. Modality-agnostic attention fusion for visual search with text feedback. *arXiv preprint arXiv:2007.00145*, 2020.
- Samuel Dooley, George Z Wei, Tom Goldstein, and John Dickerson. Robustness disparities in face detection. *Advances in Neural Information Processing Systems*, 35:38245–38259, 2022.
- Sonam Goenka, Zhaoheng Zheng, Ayush Jaiswal, Rakesh Chada, Yue Wu, Varsha Hedau, and Pradeep Natarajan. Fashionvlp: Vision language transformer for fashion retrieval with feedback. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14105–14115, 2022.
- Geonmo Gu, Sanghyuk Chun, Wonjae Kim, HeeJae Jun, Yoohoon Kang, and Sangdoo Yun. Compodiff: Versatile composed image retrieval with latent diffusion. *arXiv preprint arXiv:2303.11916*, 2023.
- Xiao Han, Sen He, Li Zhang, Yi-Zhe Song, and Tao Xiang. Uigr: unified interactive garment retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2220–2225, 2022a.
- Xiao Han, Licheng Yu, Xiatian Zhu, Li Zhang, Yi-Zhe Song, and Tao Xiang. Fashionvil: Fashion-focused vision-and-language representation learning. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXV*, pp. 634–651. Springer, 2022b.
- Xiao Han, Xiatian Zhu, Licheng Yu, Li Zhang, Yi-Zhe Song, and Tao Xiang. Fame-vil: Multi-tasking vision-language model for heterogeneous fashion tasks. *arXiv preprint arXiv:2303.02483*, 2023.
- Xintong Han, Zuxuan Wu, Phoenix X Huang, Xiao Zhang, Menglong Zhu, Yuan Li, Yang Zhao, and Larry S Davis. Automatic spatially-aware fashion concept discovery. In *Proceedings of the IEEE international conference on computer vision*, pp. 1463–1471, 2017.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *Proceedings of the International Conference on Learning Representations*, 2019.
- Phillip Isola, Joseph J Lim, and Edward H Adelson. Discovering states and transformations in image collections. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1383–1391, 2015.
- Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2901–2910, 2017.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023.
- Seungmin Lee, Dongwan Kim, and Bohyung Han. Cosmo: Content-style modulation for image retrieval with text feedback. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 802–812, 2021.

- Matan Levy, Rami Ben-Ari, Nir Darshan, and Dani Lischinski. Data roaming and early fusion for composed image retrieval. *arXiv preprint arXiv:2303.09429*, 2023.
- Linjie Li, Zhe Gan, and Jingjing Liu. A closer look at the robustness of vision-and-language pre-trained models. *arXiv preprint arXiv:2012.08673*, 2020.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023.
- Zheyuan Liu, Cristian Rodriguez-Opazo, Damien Teney, and Stephen Gould. Image retrieval on real-life images with pre-trained vision-and-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2125–2134, 2021.
- Zixian Ma, Jerry Hong, Mustafa Omer Gul, Mona Gandhi, Irena Gao, and Ranjay Krishna. Crepe: Can vision-language foundation models reason compositionally? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10910–10921, 2023.
- Long Mai, Hailin Jin, Zhe Lin, Chen Fang, Jonathan Brandt, and Feng Liu. Spatial-semantic image search by visual feature synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4718–4727, 2017.
- Vikramjit Mitra, Horacio Franco, Richard M Stern, Julien Van Hout, Luciana Ferrer, Martin Graciarena, Wen Wang, Dimitra Vergyri, Abeer Alwan, and John HL Hansen. Robust features in deep-learning-based speech recognition. *New Era for Robust Speech Recognition: Exploiting Deep Learning*, pp. 187–217, 2017.
- Andrei Neculai, Yanbei Chen, and Zeynep Akata. Probabilistic compositional embeddings for multimodal image retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4547–4557, 2022.
- Devi Parikh and Kristen Grauman. Relative attributes. In *2011 International Conference on Computer Vision*, pp. 503–510. IEEE, 2011.
- Sayak Paul and Pin-Yu Chen. Vision transformers are robust learners. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 2071–2081, 2022.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021.
- Barbara Rychalska, Dominika Basaj, Alicja Gosiewska, and Przemysław Biecek. Models in the wild: On corruption robustness of neural nlp systems. In *Neural Information Processing: 26th International Conference, ICONIP 2019, Sydney, NSW, Australia, December 12–15, 2019, Proceedings, Part III 26*, pp. 235–247. Springer, 2019.
- Kuniaki Saito, Kihyuk Sohn, Xiang Zhang, Chun-Liang Li, Chen-Yu Lee, Kate Saenko, and Tomas Pfister. Pic2word: Mapping pictures to words for zero-shot composed image retrieval. *arXiv preprint arXiv:2302.03084*, 2023.
- Soumya Sanyal, Zeyi Liao, and Xiang Ren. Robustlr: A diagnostic benchmark for evaluating logical robustness of deductive reasoners. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 9614–9631, 2022.

- Ravi Shekhar, Sandro Pezzelle, Yauhen Klimovich, Aurélie Herbelot, Moin Nabi, Enver Sangineto, and Raffaella Bernardi. Foil it! find one mismatch between image and language caption. *arXiv preprint arXiv:1705.01359*, 2017.
- Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. A corpus for reasoning about natural language grounded in photographs. *arXiv preprint arXiv:1811.00491*, 2018.
- Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. Winoground: Probing vision and language models for visio-linguistic compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5238–5248, 2022.
- Nam Vo, Lu Jiang, Chen Sun, Kevin Murphy, Li-Jia Li, Li Fei-Fei, and James Hays. 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.
- Boxin Wang, Chejian Xu, Shuohang Wang, Zhe Gan, Yu Cheng, Jianfeng Gao, Ahmed Hassan Awadallah, and Bo Li. Adversarial glue: A multi-task benchmark for robustness evaluation of language models. *arXiv preprint arXiv:2111.02840*, 2021a.
- Xiao Wang, Qin Liu, Tao Gui, Qi Zhang, Yicheng Zou, Xin Zhou, Jiacheng Ye, Yongxin Zhang, Rui Zheng, Zexiong Pang, et al. Textflint: Unified multilingual robustness evaluation toolkit for natural language processing. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations*, pp. 347–355, 2021b.
- Xuezhi Wang, Haohan Wang, and Diyi Yang. Measure and improve robustness in nlp models: A survey. *arXiv preprint arXiv:2112.08313*, 2021c.
- Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023.
- Hui Wu, Yupeng Gao, Xiaoxiao Guo, Ziad Al-Halah, Steven Rennie, Kristen Grauman, and Rogerio Feris. The fashion iq dataset: Retrieving images by combining side information and relative natural language feedback. *CVPR*, 2021.
- Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James 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.
- Jiehang Zeng, Xiaoqing Zheng, Jianhan Xu, Linyang Li, Liping Yuan, and Xuanjing Huang. Certified robustness to text adversarial attacks by randomized [mask]. *arXiv preprint arXiv:2105.03743*, 2021.
- Aimen Zerroug, Mohit Vaishnav, Julien Colin, Sebastian Musslick, and Thomas Serre. A benchmark for compositional visual reasoning. *arXiv preprint arXiv:2206.05379*, 2022.

## A Appendix

### A Creating Benchmark Datasets

We build three benchmark datasets for this work evaluating both natural corruption (CIRR-C and FashionIQ-C) and textual understanding (CIRR-D). To evaluate natural corruption with both image and text, we introduce CIRR-C and FashionIQ-C based on the existing dataset CIRR and FashionIQ. To evaluate textual understanding including variations of numerical, attributes (colour, shape and size), object removal, background and fine-grained details, we introduce CIRR-D by categorizing and expanding CIRR with synthetic images. We provide raw images and complete code for generating all types of natural corruptions and the evaluation testbed in our code zip file. For shortcut, we provide raw image link for CIRR, FashionIQ and CIRR-D. To implement **CIRR-D** dataset, both raw images and queries in different categories (numerical,

attribute, object removal, background and fine-grained variations) are provided directly. To implement **CIRR-C** and **FashionIQ-C** dataset, research can recreate the same benchmark datasets with the following steps:

1. Download CIRR and FashionIQ raw images with our provided link.
2. Preprocess image or text with the provided code of image corruption and text corruption.
3. Apply the proposed corruptions with our testbed for downstream model evaluation.

## B Sample visualization

### B.1 CIRR-C visualization

We show the visualization samples from CIRR-C in Figure. 7. Our CIRR-C is based on the CIRR dataset and implemented with both image corruptions and text corruptions. We apply 15 standard natural image corruptions, as depicted in Figure 7 (a), and demonstrate 5 levels of severity using brightness corruption as an example in Figure 7 (b). We further visualize 7 text corruptions in Figure 7 (c). For both image and text corruption, humans can easily recognize them.

### B.2 FashionIQ-C visualization

FashionIQ-C follows the same natural corruption in both image and text as in CIRR-C. We show the visualization samples from FashionIQ-C in Figure. 8. FashionIQ-C is based on the FashionIQ dataset and implemented with both image corruptions and text corruptions. We apply 15 standard natural image corruptions, as depicted in Figure 8 (a), and demonstrate 5 levels of severity using zoom blur corruption as an example in Figure 8 (b). We further visualize 7 text corruptions in Figure 8 (c).

### B.3 Textual corruption definition

In this work, we implement 7 natural textual corruptions following Rychalska et al. (2019). The definition of the textual corruptions are as follows:

- Swap: Randomly shuffles two characters within a word.
- Qwerty: Simulates errors made while writing on a QWERTY-type keyboard. Characters are swapped for their neighbors on the keyboard
- RemoveChar: Randomly removes characters from words.
- RemoveSpace: Removes a space from text, merging two words.
- Misspelling: Misspells words appearing in the Wikipedia list of commonly misspelled English words.
- Repetition: Randomly repeat words.
- Homophone: Changes words into their homophones from the Wikipedia list of common homophones. The list contains around 500 pairs or triples of homophonic words.

Examples are shown in Figure. 7 for CIRR-C and Figure. 8 for FashionIQ-C respectively.

### B.4 CIRR-D visualization

To detect textual understanding ability, we build a CIRR-D dataset with five different types of queries containing specific instructions to probe five different abilities. The source of the CIRR-D dataset is from the original CIRR, CIRR extends caption and our generated synthetic images. The triplets from the original CIRR dataset are normally with obvious variations while the synthetic triplets are normally following the

Table 4: Details of CIRRR-D dataset. The first column is the number of images. The rest columns contain the number of triplets for five probing abilities.

	Images	Numerical	Attribute	Removal	Background	Fine-grained
Val.	2297	820	1397	233	358	4181
Extend caption	-	-	-	505	812	-
Synthetic	1245	305	700	140	-	-
Total	3542	1125	2097	878	1170	4181

Table 5: Details of modality fusion for the evaluated models in this study.  $R_i, M_t, T_i$  represent reference image, modified text and target image feature respectively.  $C$  donates the composed of the two features.  $\langle, \rangle$  represents cosine similarity.

Model	$C_{RiMt}$	$C_{RiT_i}$	$C_{MtT_i}$	distance
TIRG Vo et al. (2019)	cat+residual	-	-	$\langle C_{RiMt}, T_i \rangle$
MAAF Dodds et al. (2020)	self attn+ cross attn	-	-	$\langle C_{RiMt}, T_i \rangle$
Artemis Delmas et al. (2022)	dot product	dot product	dot product	$\langle C_{RiMt}, C_{TiMt} \rangle + \langle C_{TiMt}, M_t \rangle$
CIRPLANT Liu et al. (2021)	transformer	-	-	$\langle C_{RiMt}, T_i \rangle$
CLIP4CIR Baldrati et al. (2022)	cat + residual	-	-	$\langle C_{RiMt}, T_i \rangle$
FashionViL Han et al. (2022b)	transformer	-	-	$\langle C_{RiMt}, T_i \rangle$

same structure and only local variations. Some extended caption from the original CIRRR dataset can only supply partial difference and cannot locate the target images. Therefore, we manually remove samples and retain only those triplets where the extended caption can provide sufficient variations. In detail, we visualize numerical samples in Figure 9, which is composed of triplets from the original CIRRR dataset in Figure 9 (a) and our generated synthetic triplets in Figure.9 (b). Attribute variation visualization samples are shown in Figure 10, which is composed of triplets from the original CIRRR dataset in Figure 10 (a) and our generated synthetic triplets in Figure. 10 (b). To evaluate object removal ability, the triplet source consists of three aspects. We visualize object removal triplets from the original CIRRR dataset in Figure. 11 (a), extended caption triplets in Figure. 11 (b) and our generated synthetic triplets in Figure. 11 (c). we visualize background variations samples in Figure 12, which is composed of triplets from original CIRRR dataset in Figure 12 (a) and from extended captions in Figure.12 (b). The evaluation of fine-grained variations follows the original CIRRR dataset, whose gallery set is composed of 5 highly similar images as shown in Figure. 13.

## C More experiments result

### C.1 Fine grained subset analysis of CIRRR-C

In this section, we supplement some experimental results. As shown in Figure. 4, we visualize the recall performance on the CIRRR-C subset. Comparing the subset retrieval with the whole gallery in CIRRR-C (shown in the main paper Figure.3), we can observe that subset relative robustness (range from 0.6 to 0.9) overall is higher than the whole set (range from 0.4 to 0.8). This result suggests that a smaller gallery can lead to more stable retrieval. In essence, the overall trend aligns with retrieval on all images: CLIP4CIR consistently performs the best, while IMAGE-ONLY with CLIP embedding consistently exhibits the worst retrieval performance.

### C.2 Subcategories analysis of FashionIQ-C

For detailed results in the FashionIQ-C dataset, we report the results on the three categories, namely dress, shirt and toptee respectively, as shown in Table. 6. Overall, a similar trend is observed across the three categories, with FashionViL and CLIP4CIR consistently exhibiting the highest relative robustness. In the shirt category, overall robustness tends to be slightly higher than in the dress and toptee categories. We

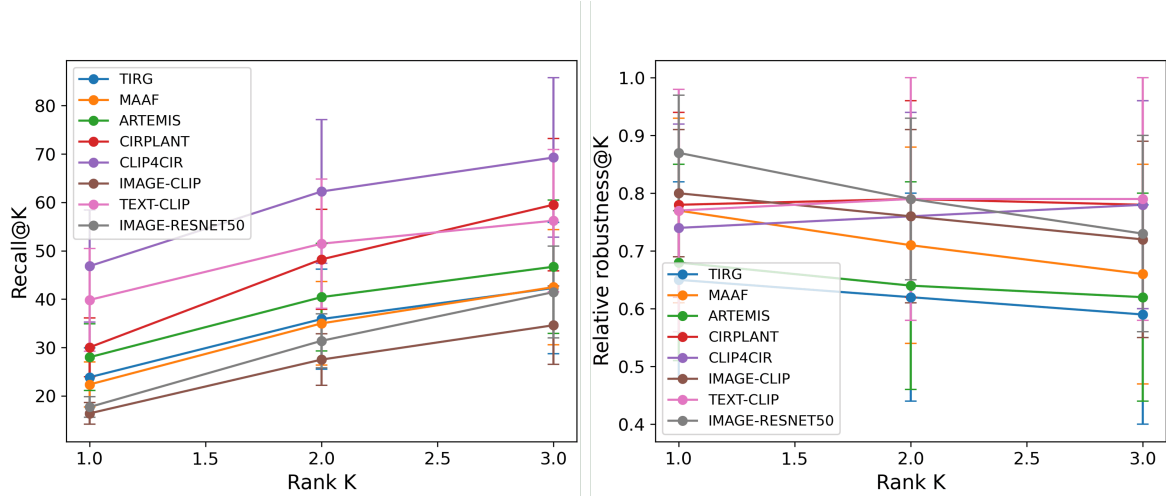


Figure 4: Models average performance in CIRR subset under 15 vision corruptions. Left: Recall vs. rank K. Right: Relative robustness vs. rank K.

Table 6: Relative robustness score for text-image composed retrieval under 15 natural image corruptions in FashionIQ-C Recall@10 for dress, shirt and toptee respectively. **Bold** is the highest relative robustness for the five composed retrieval methods.

	Noise			Blur				Weather				Digital			
<b>FashionIQ-C Dress</b>	Gauss.	Shot	Impluse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Bright	Contrast	Elastic	Pixel	JPEG
TIRG Vo et al. (2019)	0.21	0.18	0.17	0.37	0.22	0.64	0.58	0.33	0.25	0.35	0.63	0.12	0.62	0.82	0.85
MAAF Dodds et al. (2020)	0.30	0.24	0.22	0.42	0.19	0.65	0.56	0.28	0.21	0.32	0.58	0.10	0.54	0.78	0.81
ARTEMIS Delmas et al. (2022)	0.23	0.22	0.18	0.38	0.24	0.66	0.62	0.39	0.26	0.37	0.59	0.14	0.67	0.85	0.9
FashionViL Han et al. (2022b)	0.21	0.22	0.23	0.38	<b>0.34</b>	<b>0.84</b>	<b>0.72</b>	0.29	0.29	0.3	<b>0.79</b>	0.13	<b>0.88</b>	<b>1.1</b>	<b>1.1</b>
CLIP4CIR Baldrati et al. (2022)	<b>0.44</b>	<b>0.38</b>	<b>0.44</b>	<b>0.54</b>	0.24	0.74	0.52	<b>0.41</b>	<b>0.36</b>	<b>0.55</b>	0.68	<b>0.16</b>	0.42	0.75	0.82
	Noise			Blur				Weather				Digital			
<b>FashionIQ-C Shirt</b>	Gauss.	Shot	Impluse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Bright	Contrast	Elastic	Pixel	JPEG
TIRG Vo et al. (2019)	0.33	0.32	0.27	0.28	0.20	0.57	0.54	0.32	0.28	0.37	0.51	0.15	0.60	0.86	0.81
MAAF Dodds et al. (2020)	0.33	0.30	0.27	0.46	0.20	0.67	0.50	0.30	0.27	0.34	0.47	0.16	0.57	0.84	0.79
ARTEMIS Delmas et al. (2022)	0.27	0.28	0.25	0.39	<b>0.26</b>	0.62	<b>0.61</b>	0.36	0.24	0.38	0.54	0.16	0.61	0.84	0.88
FashionViL Han et al. (2022b)	0.29	0.34	0.26	0.38	<b>0.26</b>	<b>0.77</b>	0.6	0.33	0.32	0.37	0.63	0.17	<b>0.83</b>	<b>1.09</b>	<b>1.02</b>
CLIP4CIR Baldrati et al. (2022)	<b>0.47</b>	<b>0.48</b>	<b>0.45</b>	<b>0.50</b>	0.18	0.65	0.48	<b>0.51</b>	<b>0.50</b>	<b>0.65</b>	<b>0.71</b>	<b>0.27</b>	0.31	0.69	0.82
	Noise			Blur				Weather				Digital			
<b>FashionIQ-C Toptee</b>	Gauss.	Shot	Impluse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Bright	Contrast	Elastic	Pixel	JPEG
TIRG Vo et al. (2019)	0.30	0.28	0.25	0.36	0.24	0.63	0.58	0.32	0.27	0.39	0.58	0.10	0.69	0.88	0.88
MAAF Dodds et al. (2020)	0.30	0.28	0.27	0.45	0.24	0.71	0.52	0.28	0.23	0.28	0.56	0.14	0.52	0.88	0.88
ARTEMIS Delmas et al. (2022)	0.21	0.23	0.18	0.37	0.28	0.68	0.57	0.33	0.25	0.38	0.53	0.13	0.66	0.88	0.82
FashionViL Han et al. (2022b)	0.28	0.28	0.27	0.44	<b>0.32</b>	<b>0.85</b>	<b>0.69</b>	0.38	0.33	0.36	0.69	0.15	<b>0.88</b>	<b>1.09</b>	<b>1.06</b>
CLIP4CIR Baldrati et al. (2022)	<b>0.42</b>	<b>0.4</b>	<b>0.42</b>	<b>0.58</b>	0.21	0.76	0.49	<b>0.46</b>	<b>0.44</b>	<b>0.60</b>	<b>0.71</b>	<b>0.24</b>	0.39	0.78	0.84

further report the recall@10 performance on FashionIQ-C dataset as shown in Table. 7. By comparing the relative robustness in Table. 6 and corresponding recall performance in Table. 7, we can find out higher robustness doesn't mean higher recall performance. As according to the definition of relative robustness:  $\gamma = 1 - (R_c - R_p) / R_c$  following Hendrycks & Dietterich (2019), lower recall performance under clean condition  $R_c$  will lead to higher relative robustness  $\gamma$ .

Table 7: Recall@10 for text-image composed retrieval under 15 natural image corruptions in FashionIQ-C.

	Noise				Blur				Weather				Digital			
<b>FashionIQ-C</b>	Clean	Gauss.	Shot	Impluse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Bright	Contrast	Elastic	Pixel	JPEG
TIRG Vo et al. (2019)	23.8	6.6	6.1	5.4	8.1	5.3	14.6	13.5	7.7	6.3	8.8	13.8	3.0	15.2	20.3	20.2
MAAF Dodds et al. (2020)	23.4	7.2	6.4	5.9	10.4	5.0	15.8	12.3	6.6	5.5	7.3	12.7	3.1	12.7	19.4	19.4
ARTEMIS Delmas et al. (2022)	24.9	5.8	6.0	4.9	9.4	6.5	16.4	14.9	9.0	6.2	9.4	13.9	3.5	16.2	21.4	21.5
FashionViL Han et al. (2022b)	23.4	6.1	6.5	5.9	9.5	7.2	19.3	15.8	7.8	7.3	8.0	16.5	3.5	20.3	25.7	24.9
CLIP4CIR Baldrati et al. (2022)	35.9	15.9	15.2	15.6	19.4	7.5	25.7	17.8	16.5	15.6	21.5	25.2	8.2	13.3	26.5	29.7

Table 8: Relative robustness score for text-image composed retrieval under 15 natural image corruptions in COCO-C Recall@10.

	Noise			Blur				Weather				Digital			
<b>COCO-C</b>	Gauss.	Shot	Impluse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Bright	Contrast	Elastic	Pixel	JPEG
TIRG Vo et al. (2019)	0.19	0.21	0.14	0.42	0.25	0.62	0.58	0.35	0.21	0.51	0.89	0.05	0.40	0.40	0.72
ARTEMIS Delmas et al. (2022)	0.14	0.16	0.08	0.43	0.22	0.72	0.52	0.41	0.32	0.45	1.06	0.05	0.40	0.48	0.70
CLIP4CIR Baldrati et al. (2022)	0.52	0.58	0.52	0.65	0.12	0.85	0.36	0.51	0.49	0.71	0.90	0.10	0.24	0.77	0.77

### C.3 Analysis of COCO with image corruptions

To evaluate the compared models on more general domain, we implement our image corruptions on the validation set of COCO Lin et al. (2014), represented by CIRR-C. We set masked bounding box as the reference image, the raw image as the target image, and the labels of objects as modified text the following Neculai et al. (2022); Saito et al. (2023). The three compared models are trained on the CIRR dataset and evaluated on the validation set of COCO with 5000 images. The results show that large pretrained model CLIP4CIR have higher robustness than smaller models TIRG and ARTEMIS, which follow the same conclusion in paper Section 4.1.

## D Limitation

We discuss the limitation of the proposed benchmarks in this section. For benchmarking natural corruption in CIRR-C and FashionIQ-C, the method of simulating real-world corruption with the noise still has limitations. For benchmarking textual understanding in CIRR-D, it has long-tail distribution. As shown in Figure. 5, both numerical and attribute evaluation set follows the long-tail distribution. The numerical set has a large number of samples in the range of 1 to 3, while each category from 4 to 10 has only a small number of samples. The attribute evaluation set has a large number of samples with colour variations and a small number of samples with size variations. The imbalanced distribution can lead to bias towards the categories with more data. Further, we visualise the performance of the query with number one to three, four to ten respectively shown in Figure 6 left. The average recall@5 of five evaluated methods are 43.06% on number one to three, 42.36% on numbers four to ten and 44.2% on number one to ten respectively. (A sentence with multiple numbers will be categorized to multiple categories, thus number one to three and number four to ten can overlap.) Based on this subtle accuracy change, we speculate that the model also possesses a similar capability for recognizing the less frequent samples (number four to ten) in the long-tail distribution as it does for the more frequent samples (number one to three).

## E License

All the models in this study are available to the public. The model code for TIRG Vo et al. (2019) and MAAF Dodds et al. (2020) have the Apache License Version 2.0, ARTEMIS Delmas et al. (2022) has CC BY-NC-SA 4.0 License, CIRPLANT Liu et al. (2021) has MIT license and FashionViL Han et al. (2022b) has BSD License. We will provide CIRR-C, FashionIQ-C and CIRR-D publicly. These datasets are based on

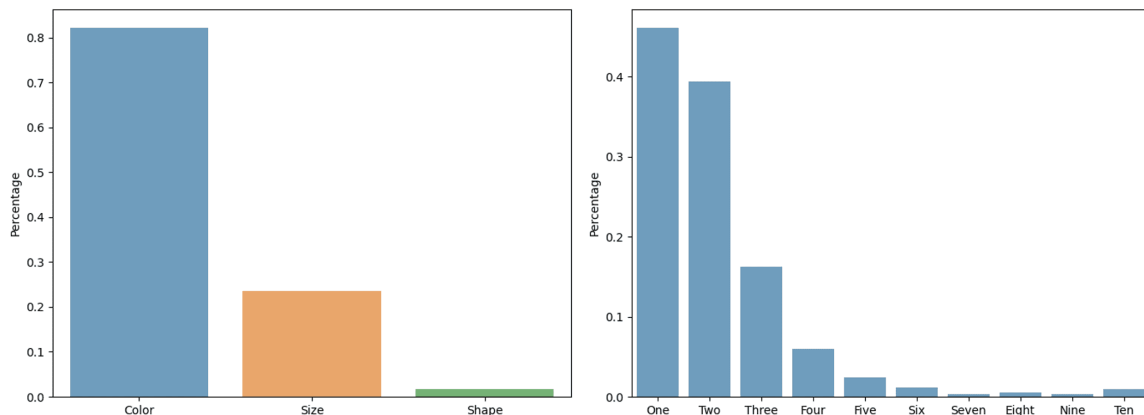


Figure 5: CIRR-D distribution for attribute variants and numerical variations.

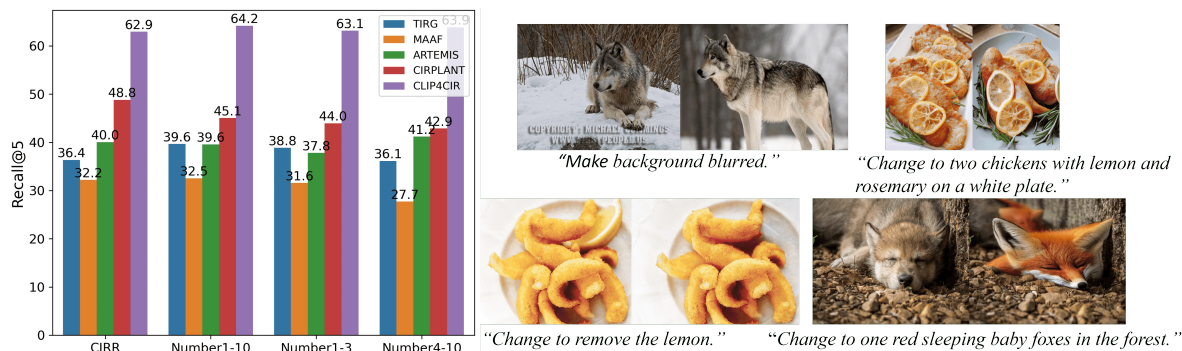
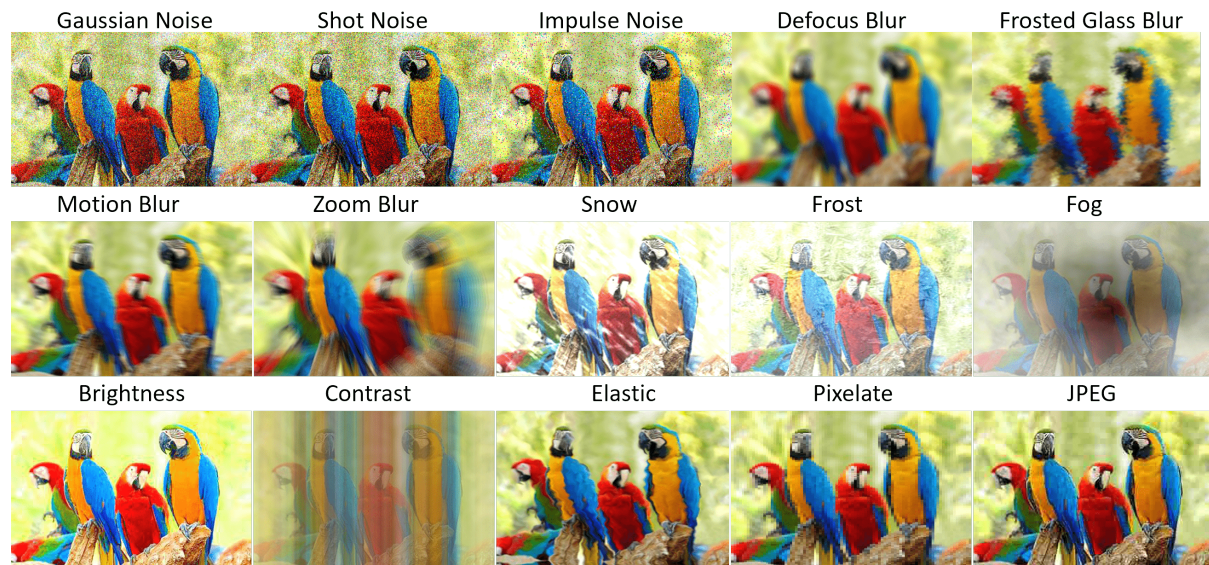


Figure 6: Left: Recall@5 on CIRR-D numerical queries. Right: Samples visualization of CIRR-D. The left and right images are the reference and target images. Except the upper left triplet, rest target images are synthetic.

existing CIRR Liu et al. (2021) and FashionIQ Wu et al. (2021). For CIRR-C and FashionIQ-C, we didn't add any new images or text sources. For CIRR-D, we further generate synthetic images and text to expand the original CIRR dataset. All of these datasets are available to the public and we apply similar licenses to our testbed code and our proposed benchmarks.





(a) Sample visualization with 15 standard image corruptions.



(b) Sample visualization of brightness corruption with 5 severities.

**Original text:**

*There were two adult dogs on the road - there was one grown puppy in the yard*

**character\_filter**

'There were two *adutl* dogs on *teh* road - there wsa oen grown puppy in the yard.'

**qwerty\_filter**

'There were two adult dogs on the road - there was one *grow5* puppy in the yard.'

**RemoveChar\_filter**

'Thre were two adult dogs on the road - *tere ws ne* grown puppy in the yard.'

**remove\_space\_filter**

'There were two adult dogs *onthe* road - there was one grown puppy in the yard.'

**misspelling\_filter**

'There *were were* two adult dogs on the road - there was one grown puppy in the yard.'

**repetition\_filter**

'There were two *adult adult* dogs on the road - there was one *grown grown* puppy in the yard.'

**homophones\_filter**

'*They're* were two adult dogs on the *rowed* - their was won grown puppy inn the yard.'

(c) Sample visualization of 7 text corruptions.

Figure 7: CIRR-C sample visualization: (a) 15 standard image corruptions, (b) 5 severities of brightness corruption and (c) 7 text corruption.



(a) Sample visualization with 15 standard image corruptions.



(b) Sample visualization of zoom blur corruption with 5 severities.

**Original text:*****Has black and white pattern.***

- character\_filter  
'has **bleka** and white pattern.'

- qwerty\_filter  
'has black and white **lattern**.'

- RemoveChar\_filter  
'has black **ad** white **attern**.'

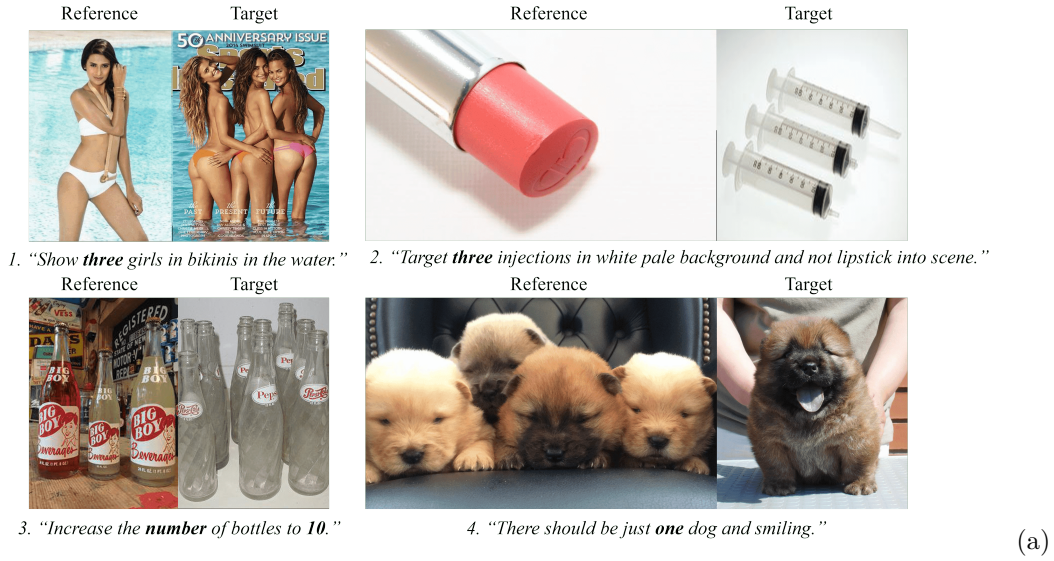
- remove\_space\_filter  
'has black **andwhite** pattern.'

- misspelling\_filter  
'has black and **white white** pattern.'

- repetition\_filter  
'has **black black black** and white pattern.'

- 20 homophones\_filter  
'has black and **wight** pattern.'





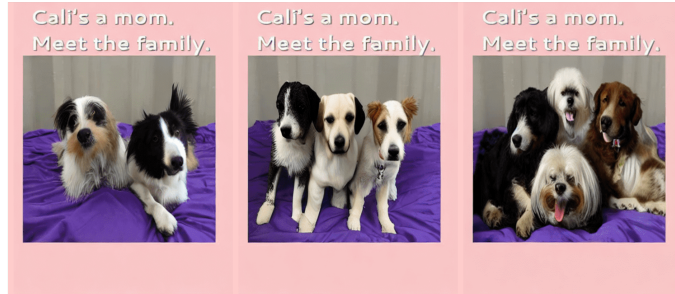
Numerical queries from original CIRRD dataset, four triplets are included.

Reference Image

Target Images



1. "Change to **two** blue and red Pepsi Colas are on the bus."
2. "Change to **three** blue and red Pepsi Colas are on the bus."
3. "Change to **four** blue and red Pepsi Colas are on the bus."



4. "Change to **two** dogs with pink backgrounds are lounging on the couch."
5. "Change to **three** dogs with pink backgrounds are lounging on the couch."
6. "Change to **four** dogs with pink backgrounds are lounging on the couch."



7. "Change to **three** sleeping baby foxes in the forest."
8. "Change to **four** sleeping baby foxes in the forest."

(b) Our generated numerical queries, eight triplets are included.

Figure 9: CIRRD sample visualization for numerical queries.



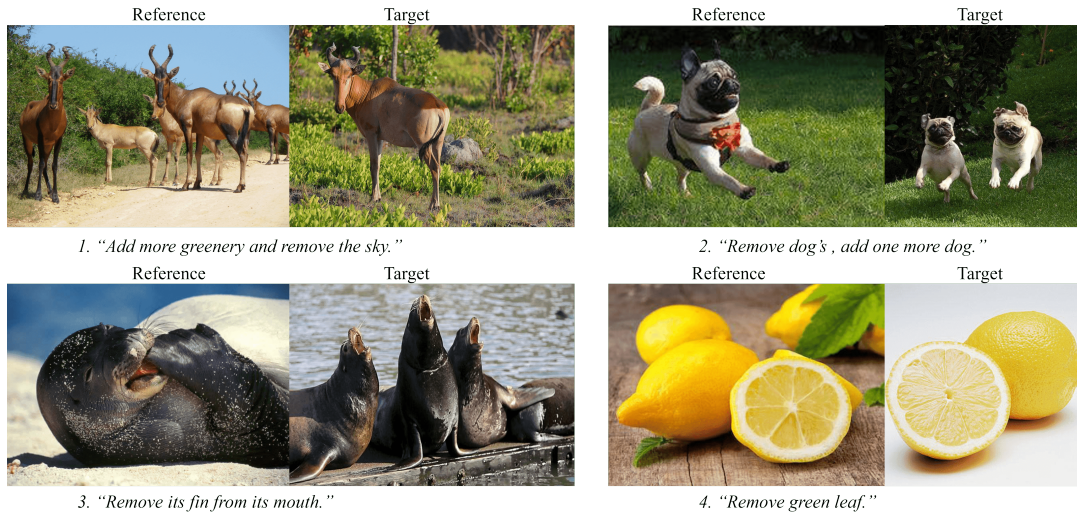
(b) Our generated attribute queries, 10 triplets are included.

Figure 10: CIRRD sample visualization for attribute queries including color, shape and size.





Object removal queries from original CIRRD dataset, four triplets are included.



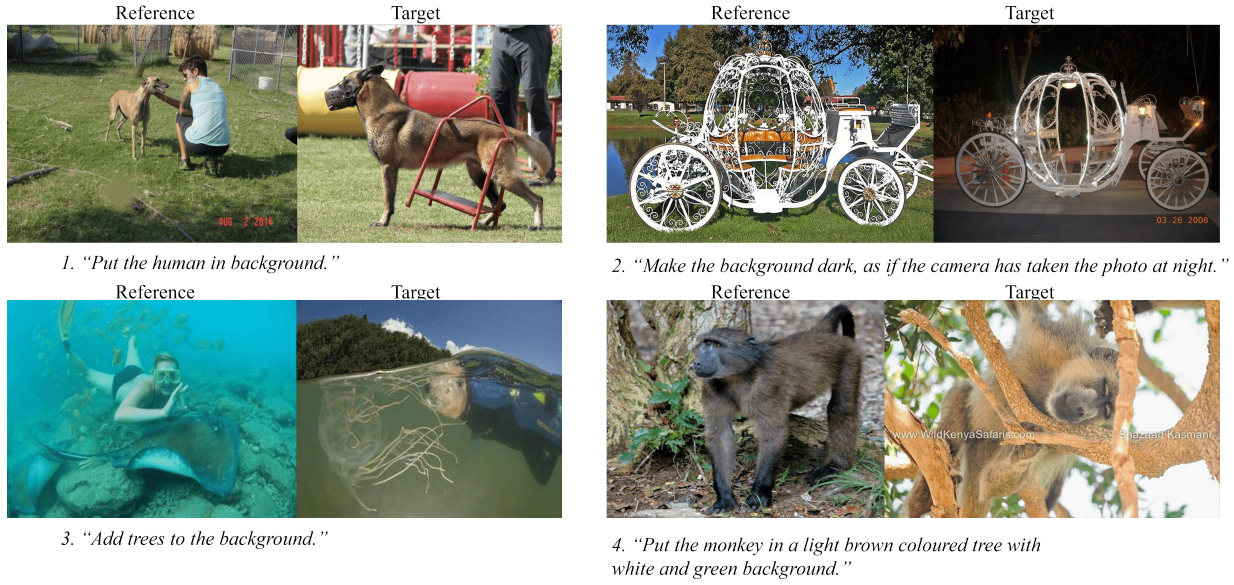
(b) Object removal queries from extend captions of original CIRRD dataset, four triplets are included.



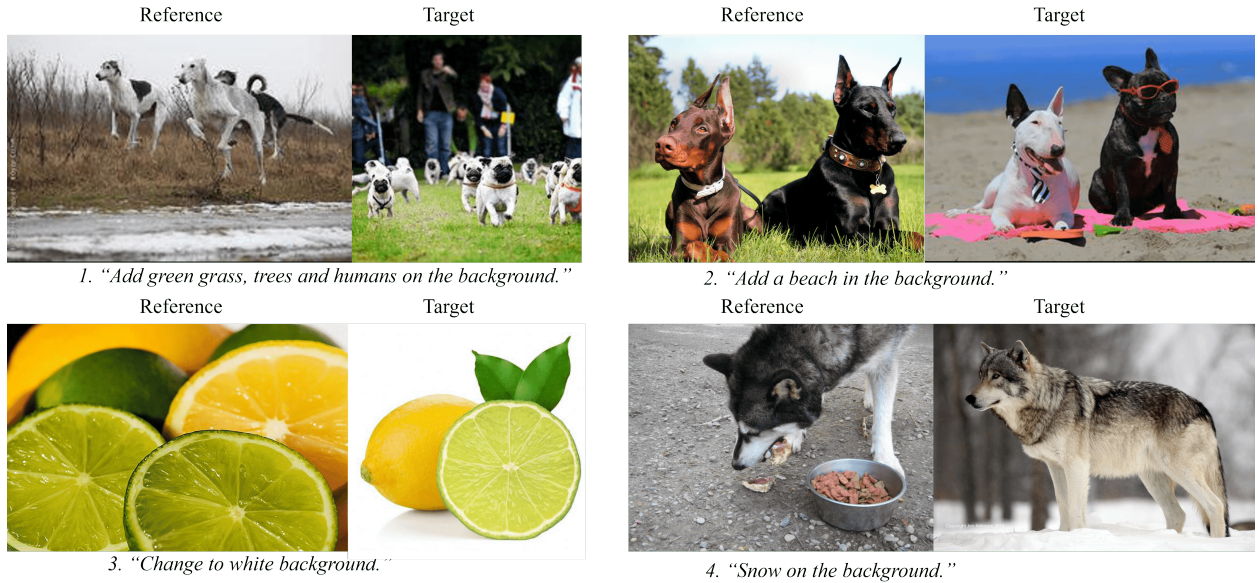
(c) Our generated object removal queries, four triplets are included.

Figure 11: CIRRD sample visualization for object removal queries.





(a) Background queries from original CIRRD dataset, four triplets are included.



(b) Background variation queries from extend captions of original CIRRD dataset, four triplets are included.

Figure 12: CIRRD-D sample visualization for background variations.

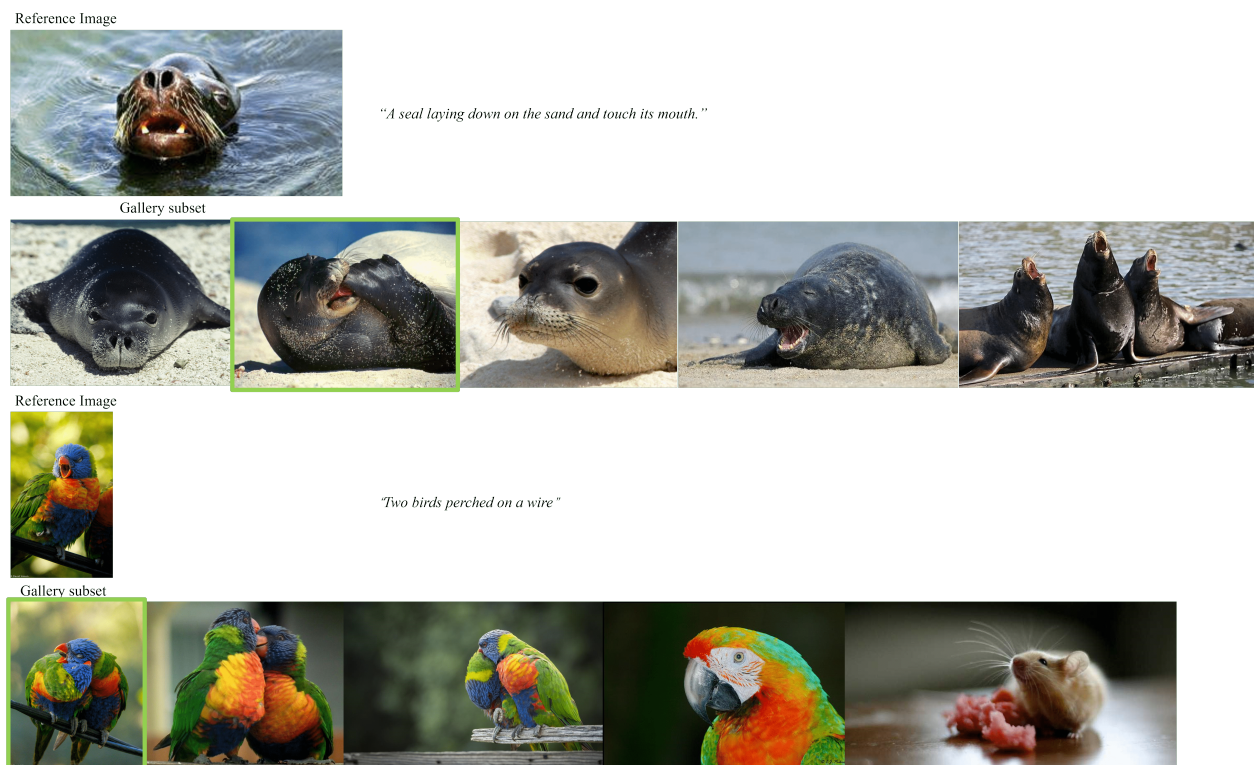


Figure 13: CIRRD sample visualization for fine-grained variation queries, 2 triplets are included. The images with a green border are the correct targets, while the other images are highly similar composing the gallery set.