

good4cir: Generating Detailed Synthetic Captions for Composed Image Retrieval

Anonymous SyntaGen submission

Paper ID 18

Abstract

001 *Composed image retrieval (CIR) enables users to search images*
 002 *using a reference image combined with textual modifications.*
 003 *Recent advances in vision-language models have improved*
 004 *CIR, but dataset limitations remain a barrier. Existing datasets*
 005 *often rely on simplistic, ambiguous, or insufficient manual*
 006 *annotations, hindering fine-grained retrieval. We introduce*
 007 *good4cir, a structured pipeline leveraging vision-language*
 008 *models to generate high-quality synthetic annotations. Our*
 009 *method involves: (1) extracting fine-grained object descriptions*
 010 *from query images, (2) generating comparable descriptions*
 011 *for target images, and (3) synthesizing textual instructions*
 012 *capturing meaningful transformations between images. This*
 013 *reduces hallucination, enhances modification diversity, and*
 014 *ensures object-level consistency. Applying our method improves*
 015 *existing datasets and enables creating new datasets across*
 016 *diverse domains. Results demonstrate improved retrieval*
 017 *accuracy for CIR models trained on our pipeline-generated*
 018 *datasets. We release our dataset construction framework to*
 019 *support further research in CIR and multi-modal retrieval.*

020 1. Introduction

021 Composed Image Retrieval (CIR) is an emerging task in vision-
 022 language research that allows users to refine image searches
 023 by providing both a reference image and a textual modification.
 024 While CIR has benefited from advancements in vision-language
 025 models (VLMs), the progress of retrieval models remains con-
 026 strained by limitations in existing datasets. Most CIR datasets
 027 are constructed through either manual annotation or automated
 028 data mining. Manually labeled datasets, such as CIRRR, provide
 029 high-quality human descriptions of modifications but are often
 030 limited in scale, expensive to create, and prone to inconsistencies
 031 in textual annotations. Automatically generated datasets, such as
 032 those based on image synthesis or retrieval-based mining, offer
 033 scalability but frequently introduce issues such as annotation
 034 noise, hallucinated content, or overly simplistic modifications
 035 that fail to capture the complexity of real-world retrieval tasks.

036 In this paper, we introduce a structured framework for
 037 generating synthetic text annotations for CIR datasets using

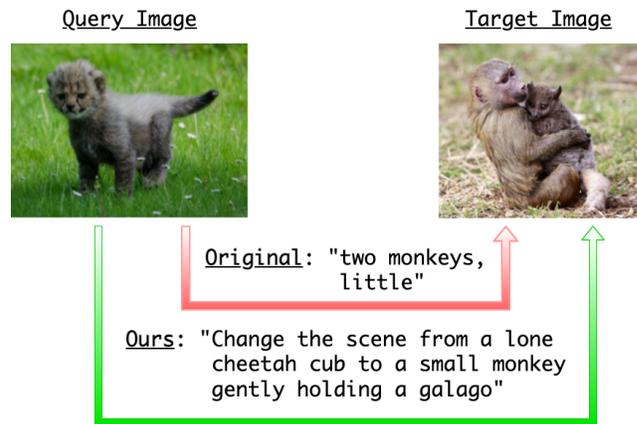


Figure 1. Existing composed image retrieval datasets are costly to construct and often have low quality text annotations. We propose a new approach that leverages VLMs to generate higher quality, synthetic text annotations for composed image retrieval.

a vision-language model-driven pipeline. Our approach
 consists of three key stages: (1) extracting detailed object-level
 descriptions from query images, (2) generating a corresponding
 set of descriptions for target images while ensuring consistency
 and capturing meaningful differences, and (3) synthesizing
 natural language modifications that describe the transformations
 required to reach the target image. This structured approach
 mitigates common pitfalls in CIR dataset construction, such
 as hallucinated object descriptions, vague or redundant
 modifications, and inconsistencies in annotation quality.

We apply our methodology to enhance existing CIR datasets
 and construct new ones across multiple domains. By evaluating
 retrieval models trained on datasets generated with our
 framework, we demonstrate improvements in retrieval accuracy,
 particularly for fine-grained modifications that require precise
 object-level reasoning. Our contributions include not only a
 scalable and effective dataset generation framework but also
 insights into the impact of dataset composition on CIR model
 performance. A GitHub link to use our dataset generation
 pipeline, to access our introduced datasets, and to re-produce
 our evaluations will be shared in our camera ready submission.

059 **2. Related Work**

060 **2.1. CIR Methods**

061 Modern composed image retrieval (CIR) methods fuse query
062 image and text representations using multimodal vision-
063 language models to retrieve relevant images [4, 5, 9, 20, 27, 31].
064 Much of the recent work focuses on algorithmic developments
065 to improve CIR performance including through the implementa-
066 tion of attention-based mechanisms [7, 36], denoising [14],
067 and interpolation-based fusion [15]. Generative vision-language
068 models [8, 19] enable training-free CIR, including video-based
069 approaches [2, 28, 30]. Textual inversion techniques [3, 13, 24]
070 learn pseudowords for query images, while other methods
071 refine cross-modal alignments [17, 25, 32, 33] for fine-grained
072 retrieval, particularly in fashion domains.

073 **2.2. CIR Datasets**

074 This paper focuses not on algorithmic developments for
075 composed image retrieval (CIR), but on CIR datasets and
076 methods for improving or creating them.

077 CIR datasets fall into two categories: manually and
078 automatically generated. Manually generated datasets include
079 CIRR [20], derived from NLVR2, which provides human
080 annotations describing image modifications. Although a key
081 benchmark, CIRR has limitations: dependence on NLVR2
082 image pairs, misaligned captions, and annotations describing
083 only single-object changes [3]. CIRCO [4] addresses these
084 issues by allowing multiple modifications per annotation,
085 sourced from MS-COCO [18], but lacks a training set and
086 serves solely for evaluation.

087 Automatically generated datasets overcome some of these
088 limitations, leveraging existing labeled data or image-generation
089 tools. Examples include LaSCo [16], synthesizing annotations
090 from large-scale datasets like VQA2.0 [12], and SynthTriplets18M [14],
091 generating images via InstructPix2Pix [6]. Domain-specific datasets,
092 such as Birds-to-Words [11] for bird species retrieval and PatternCom [22]
093 for remote sensing, also exist, alongside video retrieval datasets
094 extending CIR temporally [29, 30].

096 Most relevant to our work is MagicLens [36], which
097 constructs a dataset of 36.7 million triplets using image
098 pairs mined from web pages. After filtering duplicates and
099 low-quality content, captions and instructions are generated via
100 large multimodal and language models. While this methodology
101 is sound and the dataset could be potentially impactful for other
102 researchers working on composed image retrieval, as of March
103 2025, the dataset is not shared publicly and no code has been
104 shared to replicate it, with the authors stating on GitHub, “We
105 personally would like to release the data but the legal review
106 inside may take years.” [1]

107 Across the CIR datasets that are publicly available, there are
108 a variety of problems, regardless of the method of generation,
109 including queries where the text on its own is sufficient to find

Query Image	Target Image	Text Difference	Issue
		“show three bottles of soft drink” [20]	Query photo is unnecessary
		“has two children instead of cats” [3]	Images are not visually similar
		“Have the person be a dog” [14]	Images are too visually similar
		“Add a red ball” [4]	Modification is very simple

Figure 2. Qualitative issues with existing CIR datasets.

the target image and issues with the degree of image similarity 110
in the queries. Across existing datasets, the modifications are 111
often overly simple, focusing on a single change to a foreground 112
object. We show examples of these issues in Figure 2. Further, 113
many of the existing CIR datasets such as CIRR and CIRCO 114
are highly general in nature, lacking the specificity required 115
for many domain-specific tasks, such as medical imaging and 116
environmental monitoring. Finally, the scale of many of these 117
datasets is relatively small for any substantial training efforts. 118

3. Method 119

To improve existing CIR datasets and support the creation of 120
new ones with realistically complex textual modifications, we 121
propose good4cir, a novel pipeline that utilizes a large language 122
model – specifically OpenAI’s GPT-4o – to generate CIR 123
triplets. Our approach assumes the presence of a collection 124
of related images, which may originate from an existing 125
CIR dataset with suboptimal annotations or a novel domain 126
containing image pairs (further discussed in Section 3.6). To 127
enhance precision and reduce hallucination, we break down the 128
CIR triplet generation process into focused sub-tasks, designed 129
to encourage the production of fine-grained descriptors [10]. 130

Figure 3 depicts the structure of the proposed synthetic 131
data generation pipeline. good4cir is split into three stages, 132
which we discuss below. In the sections below, we describe the 133
general prompts for each stage. In specific domains, it may be 134
helpful to add additional specification to the prompt, such as the 135
domain of the imagery or type of scene, or a list of objects for 136
the VLM to annotate. We discuss one such case in Section 3.6, 137

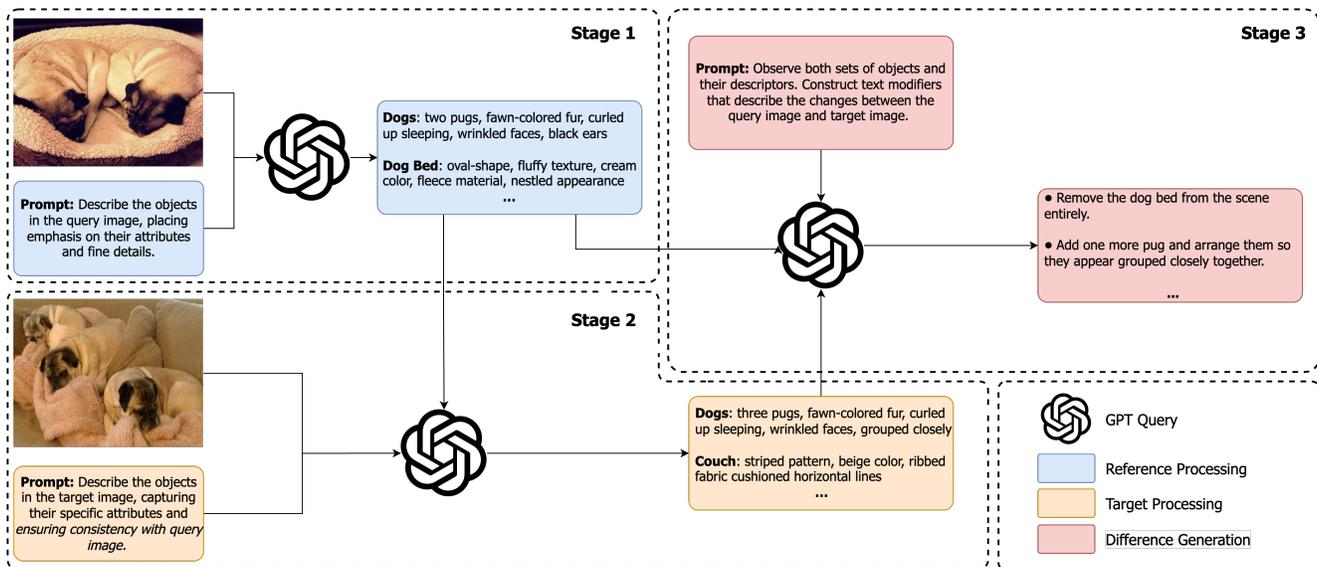


Figure 3. Our synthetic CIR data generation pipeline. The three-stage pipeline uses a structured flow of data to compare a query image and a target image without overwhelming the context window of the VLM to mitigate hallucination. In this figure, the prompts are simplified. The full prompts are discussed in the text.

138 and include the exact dataset specific prompts in the Appendix.
 139 Additionally, in Section 3.5, we demonstrate that this phased
 140 approach yields superior CIR triplets when compared with an
 141 alternative simpler approach of simply prompting a VLM to
 142 describe differences between a pair of images.

143 **3.1. Stage 1: Query Image Object Descriptions**

144 In the first stage, the VLM is prompted to generate a list of key
 145 objects and descriptors from the query image. Objects are the
 146 building blocks of any visually dense image, inherently making
 147 them signals of change. Queries used in composed image
 148 retrieval reference a specific object and a modifying caption
 149 (e.g., “Find a similar image but change the color of the chair
 150 to red”). By directing the VLM to focus on individual objects,
 151 we facilitate a more structured and detailed understanding of
 152 image differences.

153 The general form of the prompt for this stage is:

154 “Curate a list of up to X objects in the image from most
 155 prominent to least prominent. For each object, generate a
 156 list of descriptors. The descriptors should describe the exact
 157 appearance of the object, mentioning any fine-grained details.

158 Example: Object Name: [“object description 1”, “object
 159 description 2”, ..., “object description N ”]

160 Format objects and descriptors as a JSON output.”

161 The example should be constructed for the specific domain,
 162 and the quantity for X can be modified depending on the
 163 density of objects in the dataset and desired number of outputs.

164 **3.2. Stage 2: Target Image Object Descriptions**

165 In the second stage, the VLM is prompted to derive a similar list
 166 from the target image by comparing it against the list of objects
 167 from the query image, ensuring consistency and making mod-
 168 ifications when necessary. This is done by passing both the fol-
 169 lowing prompt and the output from the first stage into the VLM:

170 “Here is an image and a list of descriptors that describe a
 171 different image. Curate a similar list for this image by doing
 172 the following:

- 173 1. If there is a new object in this image that isn’t described
 174 in the description of the other image, generate a new set
 175 of descriptors.
- 176 2. If the description of an object from the other image
 177 matches the appearance of an object in this image, use
 178 the exact same list of descriptors.
- 179 3. If the object appears different in this image in comparison
 180 to the description from the other image, generate a new
 181 set of descriptors.

182 Format objects and descriptors as a JSON output.”

183 **3.3. Stage 3: Describing Differences**

184 In the final stage, the text outputs from the first two stages are
 185 passed into the VLM with the following prompt:

186 “The following are two sets of objects with descriptors that
 187 describe two different images that have been determined to
 188 be different in some ways. Analyze both lists and generate

<p>Query</p> 	<p>Single-Stage Approach</p> <ul style="list-style-type: none"> • Swap the orange curtains for green ones to create a fresher and more calming atmosphere. • Replace the orange bed skirt with a beige or light brown one to match a more neutral color scheme. • Change the chair to a darker wooden style with a more cushioned seat for a cozier touch. • Update the bedding with a subtle, patterned beige cover to create a more sophisticated and uniform look. • Replace the artwork for a simple decorative vase with fresh flowers to add color and vibrancy to the room. • Adjust the lighting to be softer and less yellow-toned to create a more inviting and relaxing ambiance. • Replace the tiled flooring with a lighter, more uniform shade to brighten up the space. • Introduce a small wooden nightstand beside the bed to provide a functional and aesthetically balanced layout.
<p>Target</p> 	<p>Three-Stage Approach</p> <ul style="list-style-type: none"> • Change the bed from single size to full size. • Raise the height of the bed from low to medium height. • Replace the decorative pillows with a single, rectangular, white pillow. • Change the decorative pillows to have a plush appearance with soft texture. • Replace the white blanket with a beige one. • Ensure the blanket maintains a thin layer, neat and tidy appearance, and smooth texture. • Replace the rust-colored bedskirt with a light brown bedskirt. • Change the bedskirt design from ruffled to plain. • Swap the rust-colored, dual-layer curtains for light green, single-layer curtains. • Introduce a headboard with a dark brown color, wooden material, and curved top. • Add a CRT model television, black in color and positioned on a stand. • Introduce a rectangular wooden table to support the TV. • Add decorative artwork of flowers with purple and pink colors, green leaves, arranged in a bouquet in a vase on the table. • Change the flooring color from light beige to light grey.

Figure 4. Comparing the direct single-stage prompting method for capturing differences, versus using good4cir’s three-stage approach.

189 *short and comprehensive instructions on how to modify the*
 190 *first image to look more like the second image. Be sure to*
 191 *mention what objects have been added, removed, or modified.*
 192 *Don’t mention “Image 1” and “Image 2” or any similar*
 193 *phrasing. Focus on having variety in the styles of captions*
 194 *that are generated, and make sure they mimic human-like*
 195 *syntactical structure and diction.”*

196 good4cir’s three-stage pipeline is aimed at addressing two
 197 fundamental issues that arise when working with VLMs:

- 198 1. **Hallucination:** VLMs generate captions that describe
 199 objects or attributes that are not actually present in the image.
 200 The multi-stage pipeline mitigates this by guiding the model
 201 to focus on concrete objects, rather than deriving a wholistic
 202 interpretation of the scene that may introduce imaginary
 203 objects or features.
- 204 2. **Limitations in Fine-Grained Captioning:** VLMs are
 205 proficient in generating relatively descriptive captions but
 206 may lack the granularity demanded by fine-grained retrieval
 207 tasks. A single-stage, direct captioning approach may lead
 208 to a vague or uninformative understanding of the object’s
 209 appearance. This idea motivates the three-stage procedure.

210 3.4. Stage 4: Caption Permutations

211 After running the first three stages, we have a dataset that
 212 consists of a number of image pairs and synthetically generated
 213 text captions describing specific differences between the images.
 214 In order to construct captions that contain more complex
 215 text differences, we implemented an automated procedure to
 216 combine individual captions into compound sentences.

217 For exactly two captions, we joined them by removing the
 218 period from the first caption, adding a comma and the word

’and’, and converting the second caption’s initial character to
 lowercase, resulting in a natural-sounding compound sentence.
 For combinations involving three captions, we sequentially
 combined the first two captions with commas, ensuring all
 intermediate captions began with lowercase letters, and added
 the conjunction ’and’ before the final caption. The final datasets
 include each original caption on its own, and then randomly
 sampled combinations of up to three captions, ensuring no
 caption was used more than once within compound sentences.
 Captions containing the verbs ‘maintain’ or ‘ensure’ were
 excluded, as they do not indicate actual differences between
 the query and target images.

We then utilized the CLIP tokenizer from OpenAI’s
 CLIP-ViT model (base-patch32) to validate each generated
 caption, discarding combinations exceeding the tokenizer’s
 77-token limit. Combination generation continued until either
 all available sentences were exhausted or a predefined limit of
 60 total combined sentences per image pair was reached.

237 3.5. Comparison to a Single-Stage Approach

238 An alternative to good4cir’s three-stage approach would be a
 239 single-stage approach, where the VLM is directly prompted
 240 to describe the differences between a pair of images. For
 241 comparison, we consider the following prompt:

242 *“The following are two different rooms that have been*
 243 *determined to be different in some ways. Analyze both lists*
 244 *and generate short instructions on how to modify the first*
 245 *image to look more like the second image. Don’t mention*
 246 *”room 1” and ”room 2” or any similar phrasing. One caption*
 247 *should discuss one modification that needs to be made to one*

Dataset	Train			Val			Test			Average Metrics	
	Image Pairs	CIR Triplets	Total Images	Image Pairs	CIR Triplets	Total Images (w/ Distractors)	Image Pairs	CIR Triplets	Total Images (w/ Distractors)	Avg. Prompt Tokens	Avg. Output Tokens
CIRR _R	28,225	199,350	16,939	4,184	22,620	2,297	–	–	–	1,600	670
Hotel-CIR	65,364	415,447	129,225	2,092	13,298	14,549	2,069	13,178	14,404	3,310	1,750

Table 1. Dataset Statistics

248 *element of the room. If one object has multiple modifications*
 249 *that need to be made, include each modification in a separate*
 250 *caption. Make sure to focus on having variety in the styles*
 251 *of captions that are generated, and make sure they mimic*
 252 *human-like conversational syntactical structure and diction.”*

253 Figure 4 compares the output of the single-stage, end-to-end
 254 approach with that of the good4cir pipeline. In the captions
 255 generated by direct captioning method, a modification to a
 256 chair in the room is described, but no chair exists in the target
 257 image. Similarly, the VLM incorrectly describes the addition
 258 of a nightstand in the second image, despite there being no
 259 nightstand. Both errors emphasize the hallucination issue with
 260 VLMs as well as their tendency to confuse objects and ideas
 261 when operating in an enlarged context window. Additionally, in
 262 the first set of captions, the model simply mentions the addition
 263 of a flower, whereas the second set provides details on the exact
 264 colors of the flowers and leaves, as well as their arrangement.
 265 This level of granularity is achieved through the structured
 266 pipeline, demonstrating the limitations of direct captioning.

267 3.6. Constructing New CIR Datasets

268 CIR datasets consist of triplets of query images, target images,
 269 and the text that describes the modification between the
 270 two. Many CIR datasets also include distractor images that
 271 are similar to the query, but do not necessarily match the
 272 text modification. Our proposed method for generating CIR
 273 captions assumes that the query-target image pairs already exist,
 274 as in the case of rewriting the captions for existing CIR datasets.

275 It is also possible to construct new CIR datasets by mining
 276 image pairs in existing image datasets that are visually similar
 277 but likely to contain differences. This is a property that is
 278 especially likely to be found in fine-grained domains, where
 279 there are large numbers of visually similar images from
 280 different classes. To mine CIR pairs from fine-grained domains,
 281 we use a combination of two different image representations:

- 282 1. **Learned Image Embedding:** Using either a domain-
 283 specific embedding model (i.e., one trained on a specific
 284 dataset) or a general-purpose model such as CLIP’s image
 285 encoder, we can identify the most semantically similar
 286 image for each image in a dataset. This process generates
 287 pairs of related images based on the similarity notion that
 288 was optimized over during the model training.
- 289 2. **Perceptual Hashing:** We use perceptual hashing and select
 290 both a minimum and maximum hash distance, allowing

us to identify pairs that structurally and visually similar,
 without being identical.

The exact similarity thresholds, and relative importance of
 the learned image similarity and perceptual hash similarity vary
 as a function of the dataset.

296 4. Datasets

297 We use our proposed approach to generate synthetic text
 298 annotations for two new datasets – CIRR_R, which is a re-written
 299 version of the CIRR dataset, and Hotel-CIR, a new CIR dataset
 300 focused on hotel recognition, a very object-centric fine-grained
 301 problem domain. Table 1 includes details on the number of
 302 image pairs, generated CIR triplets and total images (including
 303 distractors) in the training, validation and test sets, as well as
 304 the average number of GPT-4o tokens used per prompt.

305 4.1. CIRR_R

306 We use our approach to re-write the captions for the CIRR
 307 training and validation sets. As of March 2025, using the gpt-4o
 308 model and the OpenAI Batch API, it cost just about \$200 to
 309 generate all of the synthetic captions for CIRR_R.

310 Figure 5 (top) shows several examples of image pairs
 311 from the original CIRR dataset with the original CIRR text
 312 difference caption, and a sample of our re-written captions.
 313 These examples show that not only does our proposed approach
 314 generate many text prompts for each image query, but those
 315 prompts are also significantly richer in both the variations they
 316 describe and the language and grammar that they use to describe
 317 them. Additional examples can be found in the Appendix.

318 The CIRR test set is not publicly shared. This limits the
 319 relevance of our re-written captions for evaluating performance
 320 on the CIRR test set, as those captions are still in the same
 321 style as the original training set – however, in Section 6 we
 322 show that training on the rewritten dataset yields performance
 323 improvement on the zero-shot CIR dataset CIRCO.

324 4.2. Hotel-CIR

325 In order to construct the Hotel-CIR dataset, we start from the
 326 Hotels-50K dataset [26]. The hotels domain is ideal for this
 327 pipeline because the scenes in the images are dense in terms
 328 of the number of objects in any given image, and there are
 329 large numbers of visually similar images, requiring CIR models
 330 to learn subtle visual differences and rich representations of
 331 textual and semantic features.

Query Image	Target Image	Text Modifiers
		<p>CIRR:</p> <ul style="list-style-type: none"> • have an avocado in the background <p>CIRR_R:</p> <ul style="list-style-type: none"> • Remove the glass of chocolate smoothie with banana slice garnish • Replace the banana with an avocado sliced in half, showing light green flesh and a large seed • Change the chocolate smoothie to a yellow smoothie with a creamy texture, pale yellow color, and tiny bubbles on top • Swap the transparent glass for a mason jar with an embossed logo, open top, and add a green and white spiral-patterned straw • Add a circular sticker label with multicolored text reading 'DE LIGHT FUL' affixed near the top • Introduce a flat, light brown, smooth wooden surface with natural grain lines as the background
		<p>CIRR:</p> <ul style="list-style-type: none"> • instead of rabbits dog is sitting in grass <p>CIRR_R:</p> <ul style="list-style-type: none"> • Remove both guinea pigs from the scene • Introduce a tan and white short-haired dog with alert ears and a black nose in a reclining position • Add a bright green fern with feathery leaves behind the dog to enrich the backdrop • Include a cylindrical tree trunk with brown bark beside the dog to enhance the natural setting • Alter the vivid green grass to appear in tandem with the new elements, supporting a cohesive natural environment
		<p>Hotel-CIR:</p> <ul style="list-style-type: none"> • Remove the bathtub, Add a toilet with a white porcelain material, standard size, compact shape, attached tank, and chrome handle • Place a white rectangular bath mat with a soft texture and non-slip backing on the floor • Add decorative plants with green leaves in small white pots and place them on the countertop • Install a wooden door with light brown color, modern handle, and smooth surface • Ensure the door has a hinged design • Ensure the decorative plants are artificial and neatly arranged • Ensure the added door has a clean appearance and light finish • Remove the shower curtain
		<p>Hotel-CIR:</p> <ul style="list-style-type: none"> • Replace the tan-colored blanket with a beige-tan blanket on the bed, Remove the two pillows from the bed • Change the nightstand to have a medium brown color instead of tan • Modify the nightstand to have two drawers instead of a single drawer, • Replace the simple headboard of the bed with a headboard that has medium-height horizontal slats and a smooth texture • Add a wooden table with three drawers, matching the nightstand and bed frame, positioned against the wall • Include two pieces of floral artwork with green and beige coloring, framed and mounted on the wall side by side above the wooden table • Install two mounted lights on the walls • Ensure no bedskirt is visible around the bed • Make sure that the bedspread on the bed is centered

Figure 5. Example generated text differences for the CIRR_R (top) and Hotel-CIR (bottom) using our synthetic data generation pipeline. For CIRR_R, we include the original caption as well.

332 We construct (query, target) image pairs by first computing
 333 image embeddings for the images in the Hotels-50K dataset
 334 using the pre-trained model from [34], and selecting the nearest
 335 neighbor for each image that is not from the same hotel (to
 336 guarantee that there are possible modifications to describe in
 337 text). We then use perceptual hashing to filter image pairs that
 338 are either too dissimilar or nearly identical, using a similarity
 339 threshold between 25 and 35 (inclusive). Near identical matches
 340 can occur in the original Hotels-50K dataset, as different hotels

in the same chain occasionally use the same promotional
 images. Combining the learned image similarity and the
 perceptual hashing thresholding yields a set of image pairs that
 can be passed through the synthetic data pipeline to generate
 data triplets of a CIR dataset.

The specific prompts used at each stage of the pipeline to
 generate the Hotel-CIR dataset can be found in the Appendix.
 These captions are slightly modified from the “general” case,
 as including domain-specific information (such as the fact that

341
 342
 343
 344
 345
 346
 347
 348
 349

350 these images come from hotel rooms, and providing a list of
351 specific objects of interest) yields improved text differences.

352 We additionally include distractor images in the Hotel-CIR
353 dataset. To find reasonable distractor images, we embed the
354 entire Hotels-50K dataset using OpenAI’s CLIP-ViT image
355 encoder (base-patch32). For every (query, target) pair in the
356 proposed dataset, we find any other images that have higher
357 cosine similarity in the CLIP image embedding space than the
358 query and target. We randomly sample up to 5 of these images
359 as distractors for every image pair in CIR. The same image may
360 be a distractor for multiple pairs. These distractors ensure that
361 the composed image retrieval task in this dataset is challenging,
362 and that models trained on it must actually learn to incorporate
363 the information from the text difference caption, rather than
364 simply finding visually similar image pairs.

365 Figure 5 (bottom) shows several examples of CIR triplets
366 from this new dataset, and additional examples can be found
367 in the Appendix.

368 5. Evaluation

369 To demonstrate how effective our proposed pipeline is at
370 generating high-quality data, we conduct a series of experiments
371 training simple CIR models on both existing datasets and
372 our synthesized datasets created using the good4cir approach.
373 We train supervised models based on the CLIP [23] ViT-B
374 backbone. We train three modules: an image encoder f_I , a text
375 encoder f_T , and a multimodal fusion mechanism f_F , where
376 f_I, f_T are the CLIP image and text ViT-B models, respectively.
377 f_F is implemented using 4 sequential cross attention layers
378 using the text tokens as Q and the image tokens and previous
379 outputs as KV , followed by an attentional pooling as defined
380 by Yu et al. [35]. We define a forward pass through the entire
381 model as $f(Q, M) = f_F(f_I(Q), f_T(M))$ for a query image
382 and modification text pair Q, M . This model is optimized
383 contrastively with the following loss function, given a batch of
384 size N , $\{(Q_i, M_i, T_i), i \in \{1, 2, \dots, N\}\}$:

$$\mathcal{L} = \frac{\exp(\text{sim}(f(Q_i, M_i), f_I(T_i)) / \tau)}{\sum_{j=1}^N \exp(\text{sim}(f(Q_j, M_j), f_I(T_j)) / \tau)}$$

385 This framework is optimized with AdamW [21] with a
386 weight decay of $1e-2$.

387 We trained this model on the following datasets and their
388 combinations:

- 389 1. CIRR (baseline): Composed Image Retrieval on Real-life
390 images dataset.
- 391 2. CIRR_R: a variant of the CIRR dataset rewritten using the
392 proposed pipeline.
- 393 3. Hotel-CIR: a composed image retrieval dataset generated
394 for the hotels domain using the VLM-powered pipeline.

395 Because the good4cir pipeline generates a number of cap-
396 tions for every (query, target) image pair, the CIRR_R dataset
397 includes a significantly larger number of triplets than the original

Method	R@1	R@2	R@5	R@10	R@50
CIRR	16.506	25.205	41.181	56.289	82.072
CIRR _R	9.470	16.337	29.759	43.398	72.265
CIRR + CIRR _R	19.181	29.976	47.566	61.157	86.048

Table 2. Evaluation on CIRR test set. We evaluate CIRR_R and Hotel-CIR against CIRR (baseline) using a performance metric of Recall@K (or R@K). The best results are bolded.

CIRR dataset. To ensure fairness in our evaluation, when we
train on CIRR_R, we sample the synthetic captions and only in-
clude a single caption for each image pair. It likely would be ben-
eficial to train on the full dataset, but that would make the com-
parison between models trained on CIRR and CIRR_R unfair.

6. Results

To evaluate the quality of the data produced by the good4cir
pipeline, we compare retrieval performance across various
training setups: (1) models trained on existing CIR datasets
(CIRR), (2) models trained on good4cir generated datasets
(CIRR_R, Hotel-CIR), and (3) models trained on a combination
of both dataset types. All model setups were evaluated on the
Hotel-CIR, CIRR, and CIRCO test sets. The results from these
experiments are summarized in Tables 2, 3, and 4.

6.1. CIRR Evaluation

Table 2 summarizes the results from training on the CIRR,
CIRR_R, and their aggregate datasets and evaluating on the
original CIRR test set. Training with only CIRR_R captions
degrades retrieval performance compared to training on the
original CIRR training set. Since the text modifiers in the
CIRR_R dataset were reformulated to introduce greater semantic
complexity, they are no longer well aligned with the query
composition of the CIRR test set. Consequently, the model
struggles to align text queries to their corresponding images.
However, when CIRR and CIRR_R are combined, the model
exceeds that of the CIRR baseline, suggesting that the diverse
captioning offered by the CIRR_R strengthens the model’s
ability to generalize when integrated with CIRR.

6.2. Hotel-CIR Evaluation

The model trained only on the original CIRR captions, and eval-
uated on the Hotel-CIR test set achieves the lowest recall scores
across all thresholds, signifying its limitations in fine-grained
composed image retrieval tasks. By comparison, training on
only CIRR_R data offers a small boost in performance which
is most apparent at higher recall levels. However, the retrieval
accuracy achieved when coupling these datasets together
surpasses that of any one dataset alone. It is reasonable to
assume that the model benefits from the greater diversity in
length, complexity, and style of training examples provided by
the combined training set.

Method	R@5	R@10	R@50	R@100
CIRR	1.27	2.03	5.80	9.07
CIRR _R	1.61	2.75	7.52	11.22
CIRR + CIRR _R	2.07	3.20	8.66	13.09
Hotel-CIR	8.32	12.35	26.07	34.41
CIRR + Hotel-CIR	7.85	11.77	24.72	32.70
CIRR _R + Hotel-CIR	8.62	12.23	25.73	34.15
CIRR + CIRR _R + Hotel-CIR	8.57	12.20	25.69	34.04

Table 3. Evaluation on Hotel-CIR test set. We evaluate training on CIRR (baseline), CIRR_R and Hotel-CIR using the performance metric of Recall@K. The best results are bolded.

438 Still, training exclusively on Hotel-CIR data yields the
 439 greatest performance boost. Given that it is a domain-specific
 440 dataset that places an emphasis on small, object-level modifi-
 441 cations, Hotel-CIR better guides the model in understanding
 442 subtle visual differences. As shown in Table 3, Hotel-CIR
 443 achieves the highest recall accuracies at R@10, R@50, R@100,
 444 and third highest at R@5. This is likely due to the CIRR_R
 445 introducing specific concepts that help retrieval in a few select
 446 cases. Otherwise, coupling the Hotel-CIR dataset with any
 447 data set from the CIRR domain (CIRR or CIRR_R) negatively
 448 impacts retrieval performance. Since the concepts of the CIRR
 449 domain have minimal overlap with the hotels domain, they
 450 likely disrupt the patterns that the model is trying to learn from
 451 hotel-related images, introducing noise into the model.

452 6.3. CIRCO Evaluation

453 CIRCO is a zero shot composed image retrieval dataset that has
 454 multiple possible targets per query. In comparison to CIRR, the
 455 CIRCO captions are generally longer and more descriptive in
 456 their composition, making its test set a more relevant evaluation
 457 for the utility of the good4cir-generated datasets than the
 458 original CIRR test set.

459 Table 4 shows results on the CIRCO test set when training
 460 with CIRR, CIRR_R and their combinations, as well as
 461 combining them with the Hotel-CIR dataset for a single more
 462 expansive training dataset. Training on the CIRR_R dataset
 463 exceeds the performance of training only on CIRR, and
 464 combining them together achieves slightly better performance
 465 still. This indicates that CIRR_R is better aligned with the textual
 466 structure and complexities of the CIRCO test set than CIRR. We
 467 further demonstrate this by training on the aggregate of CIRR,
 468 CIRR_R, and Hotel-CIR which nearly doubles the mAP score at
 469 mAP@5, mAP@10, mAP@50, and mAP@100. These results
 470 suggest that the captions generated by the good4cir pipeline
 471 improve the model’s ability to generalize across different
 472 retrieval tasks of varying complexities.

473 7. Limitations

474 While the proposed approach to generating synthetic text
 475 annotations for CIR datasets mitigates known limitations of

Method	mAP@5	mAP@10	mAP@25	mAP@50
CIRR	2.54	2.78	3.14	3.54
CIRR _R	2.72	3.29	3.84	4.12
CIRR + CIRR _R	2.84	3.43	4.21	4.60
CIRR + CIRR _R + Hotel-CIR	4.64	5.39	6.38	7.04

Table 4. Evaluation on CIRCO test set. We evaluate training on CIRR (baseline), CIRR_R and Hotel-CIR using the performance metric of mAP@K. The best results are bolded.

VLMs, several challenges persist: 476

- 477 • **Hallucination:** The three-stage pipeline reduces but does
 478 not fully eliminate hallucination. Particularly when query
 479 and target images are highly similar, the VLM occasionally
 480 describes objects not present in either image. Hallucinations
 481 are less frequent in datasets with more visually distinct image
 482 pairs (e.g., CIRR dataset). 482
- 483 • **Counting:** VLMs often inaccurately count objects, resulting
 484 in captions that correctly identify objects but incorrectly
 485 specify their quantity. 485
- 486 • **Sentence Structure:** Despite prompts requesting varied
 487 styles, chat-based LLM outputs often exhibit limited stylistic
 488 diversity. Future work could address this by adding a
 489 post-processing step to rewrite captions in diverse styles. 489
- 490 • **Object-centric Focus:** The pipeline primarily captures
 491 variations in individual objects, limiting its effectiveness
 492 for non-object-centric datasets and abstract, conceptual
 493 differences. For instance, it might describe furniture changes
 494 in a room but miss broader shifts, such as from a modern to
 495 a traditional ambiance. 495
- 496 • **Cost:** The proposed method relies on OpenAI’s GPT-4o,
 497 incurring a per-query cost. While substantially cheaper than
 498 human annotation, this expense remains noteworthy. We
 499 explored open-source VLM alternatives but found GPT-4o
 500 significantly superior. 500

501 8. Conclusion

502 In this work, we presented good4cir, a structured and scalable
 503 pipeline for generating synthetic, high-quality text annotations
 504 for Composed Image Retrieval datasets. By leveraging advanced
 505 vision-language models and a carefully designed multi-stage
 506 prompting strategy, our approach generates richer and more di-
 507 verse textual annotations than existing datasets. We introduced
 508 two new datasets, CIRR and Hotel-CIR, created using good4cir,
 509 and demonstrated through evaluations on composed image re-
 510 trieval benchmarks that training with these datasets improves
 511 composed image retrieval accuracy in general. Our datasets and
 512 publicly available construction framework, which can be found
 513 at [https://github.com/tbd/after/camera/](https://github.com/tbd/after/camera/ready)
 514 [ready](https://github.com/tbd/after/camera/ready) aim to facilitate further progress and innovation in com-
 515 posed image retrieval and broader multimodal retrieval research. 515

516

References

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

- [1] <https://github.com/google-deepmind/magiclens/issues/8>
- [2] Bain, M., Nagrani, A., Varol, G., Zisserman, A.: Frozen in time: A joint video and image encoder for end-to-end retrieval. In: IEEE International Conference on Computer Vision (2021)
- [3] Baldrati, A., Agnolucci, L., Bertini, M., Del Bimbo, A.: Zero-shot composed image retrieval with textual inversion. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 15338–15347 (October 2023)
- [4] Baldrati, A., Bertini, M., Uricchio, T., Bimbo, A.D.: Composed image retrieval using contrastive learning and task-oriented clip-based features. ACM Transactions on Multimedia Computing, Communications and Applications **20**(3), 1–24 (2023)
- [5] 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>
- [6] Brooks, T., Holynski, A., Efros, A.A.: Instructpix2pix: Learning to follow image editing instructions. arXiv preprint arXiv:2211.09800 (2022)
- [7] Chen, Y., Gong, S., Bazzani, L.: Image search with text feedback by visiolinguistic attention learning. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2998–3008 (2020). <https://doi.org/10.1109/CVPR42600.2020.00307>
- [8] Dai, W., Li, J., Li, D., Tiong, A.M.H., Zhao, J., Wang, W., Li, B., Fung, P., Hoi, S.: Instructclip: towards general-purpose vision-language models with instruction tuning. In: Proceedings of the 37th International Conference on Neural Information Processing Systems. NIPS '23, Curran Associates Inc., Red Hook, NY, USA (2024)
- [9] 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 (2022)
- [10] Flemings, J., Zhang, W., Jiang, B., Takhirov, Z., Annavam, M.: Characterizing context influence and hallucination in summarization (2024), <https://arxiv.org/abs/2410.03026>
- [11] Forbes, M., Kaeser-Chen, C., Sharma, P., Belongie, S.: Neural naturalist: Generating fine-grained image comparisons. arXiv preprint arXiv:1909.04101 (2019)
- [12] Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., Parikh, D.: Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
- [13] Gu, G., Chun, S., Kim, W., Kang, Y., Yun, S.: Language-only training of zero-shot composed image retrieval. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2024)
- [14] Gu, G., Chun, S., Kim, W., Jun, H., Kang, Y., Yun, S.: Compodiff: Versatile composed image retrieval with latent diffusion. Transactions on Machine Learning Research (2024), <https://openreview.net/forum?id=mKtlzW0bWc>, expert Certification
- [15] Jang, Y.K., Huynh, D., Shah, A., Chen, W.K., Lim, S.N.: Spherical linear interpolation and text-anchoring for zero-shot composed image retrieval. arXiv preprint arXiv:2405.00571 (2024), <https://arxiv.org/abs/2405.00571>
- [16] Levy, M., Ben-Ari, R., Darshan, N., Lischinski, D.: Data roaming and quality assessment for composed image retrieval. In: AAAI. vol. 38 (2024). <https://doi.org/10.1609/aaai.v38i4.28081>, <https://ojs.aaai.org/index.php/AAAI/article/view/28081>
- [17] Li, D., Zhu, Y.: Visual-linguistic alignment and composition for image retrieval with text feedback. In: 2023 IEEE International Conference on Multimedia and Expo (ICME). pp. 108–113 (2023). <https://doi.org/10.1109/ICME55011.2023.00027>
- [18] 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: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) Computer Vision – ECCV 2014. pp. 740–755. Springer International Publishing, Cham (2014)
- [19] Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. In: Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., Levine, S. (eds.) Advances in Neural Information Processing Systems. vol. 36, pp. 34892–34916. Curran Associates, Inc. (2023)
- [20] 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 (ICCV). pp. 2125–2134 (October 2021)
- [21] Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: International Conference on Learning Representations (2017), <https://api.semanticscholar.org/CorpusID:53592270>
- [22] Psomas, B., Kakogeorgiou, I., Efthymiadis, N., Tolia, G., Chum, O., Avrithis, Y., Karantzalos, K.: Composed image retrieval for remote sensing. In: IGARSS 2024 - 2024 IEEE International Geoscience and Remote Sensing Symposium (2024)
- [23] 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

621 language supervision. In: International conference on
622 machine learning. pp. 8748–8763. PMLR (2021)

623 [24] Saito, K., Sohn, K., Zhang, X., Li, C.L., Lee, C.Y., Saenko,
624 K., Pfister, T.: Pic2word: Mapping pictures to words for
625 zero-shot composed image retrieval. CVPR (2023)

626 [25] Song, C.H., Hwang, T., Yoon, J., Choi, S., Gu, Y.H.:
627 Syncmask: Synchronized attentional masking for fashion-
628 centric vision-language pretraining. In: Proceedings of the
629 IEEE/CVF Conference on Computer Vision and Pattern
630 Recognition. pp. 13948–13957 (2024)

631 [26] Stylianou, A., Xuan, H., Shende, M., Brandt, J., Souvenir,
632 R., Pless, R.: Hotels-50k: A global hotel recogni-
633 tion dataset. In: The AAAI Conference on Artificial
634 Intelligence (AAAI) (2019)

635 [27] Sun, S., Ye, F., Gong, S.: Training-free zero-
636 shot composed image retrieval with local concept
637 reranking. arXiv preprint arXiv:2312.08924 (2024),
638 <https://arxiv.org/abs/2312.08924>

639 [28] Ventura, L., Yang, A., Schmid, C., Varol, G.: CoVR-2:
640 Automatic Data Construction for Composed Video
641 Retrieval. IEEE Transactions on Pattern Analysis
642 & Machine Intelligence **46**(12), 11409–11421 (Dec
643 2024). <https://doi.org/10.1109/TPAMI.2024.3463799>,
644 [https://doi.ieeecomputersociety.org/
645 10.1109/TPAMI.2024.3463799](https://doi.ieeecomputersociety.org/10.1109/TPAMI.2024.3463799)

646 [29] Ventura, L., Yang, A., Schmid, C., Varol, G.: CoVR-2:
647 Automatic data construction for composed video retrieval.
648 IEEE TPAMI (2024)

649 [30] Ventura, L., Yang, A., Schmid, C., Varol, G.: CoVR:
650 Learning composed video retrieval from web video
651 captions. AAAI (2024)

652 [31] Vo, N., Jiang, L., Sun, C., Murphy, K., Li, L.J., Fei-Fei, L.,
653 Hays, J.: Composing text and image for image retrieval-an
654 empirical odyssey. In: Proceedings of the IEEE/CVF
655 conference on computer vision and pattern recognition.
656 pp. 6439–6448 (2019)

657 [32] Wan, Y., Wang, W., Zou, G., Zhang, B.: Cross-modal
658 feature alignment and fusion for composed image
659 retrieval. In: Proceedings of the IEEE/CVF Conference
660 on Computer Vision and Pattern Recognition (CVPR)
661 Workshops. pp. 8384–8388 (June 2024)

662 [33] Xu, Y., Bin, Y., Wei, J., Yang, Y., Wang, G., Shen, H.T.:
663 Align and retrieve: Composition and decomposition
664 learning in image retrieval with text feedback. IEEE
665 Transactions on Multimedia **26**, 9936–9948 (2024).
666 <https://doi.org/10.1109/TMM.2024.3417694>

667 [34] Xuan, H., Stylianou, A., Pless, R.: Improved embeddings
668 with easy positive triplet mining. In: Proceedings of
669 the IEEE/CVF Winter Conference on Applications of
670 Computer Vision. pp. 2474–2482 (2020)

671 [35] Yu, J., Wang, Z., Vasudevan, V., Yeung, L., Seyedhosseini,
672 M., Wu, Y.: Coca: Contrastive captioners are image-text
673 foundation models (2022)

[36] Zhang, K., Luan, Y., Hu, H., Lee, K., Qiao, S., Chen, W.,
Su, Y., Chang, M.W.: MagicLens: Self-supervised image
retrieval with open-ended instructions. In: Salakhutdinov,
R., Kolter, Z., Heller, K., Weller, A., Oliver, N., Scarlett,
J., Berkenkamp, F. (eds.) PMLR. Proceedings of
Machine Learning Research, vol. 235 (21–27 Jul 2024),
[https://proceedings.mlr.press/v235/
zhang24an.html](https://proceedings.mlr.press/v235/zhang24an.html)