

# Benchmarking Composed Image Retrieval for Applied Earth Observation

Bill Psomas<sup>a</sup>, Dionysis Christopoulos<sup>b</sup>, Thanasis Petropoulos<sup>b</sup>, Nikos Efthymiadis<sup>a</sup>, Ioannis Kakogeorgiou<sup>c</sup>, Ondřej Chum<sup>a</sup>, Yannis Avrithis<sup>d</sup>, Giorgos Toliás<sup>a</sup> and Konstantinos Karantzas<sup>b</sup>

<sup>a</sup>Visual Recognition Group, Department of Cybernetics, Czech Technical University in Prague, Czechia

<sup>b</sup>Remote Sensing Laboratory, School of Rural, Surveying and Geoinformatics Engineering, National Technical University of Athens, Greece

<sup>c</sup>Institute of Informatics & Telecommunications, National Centre for Scientific Research “Demokritos”, Greece

<sup>d</sup>Department of Informatics and Telecommunications, National and Kapodistrian University of Athens, Greece

## ARTICLE INFO

### Keywords:

Remote Sensing Image Retrieval  
Composed Image Retrieval  
Multimodal Retrieval  
Vision–Language Models  
Earth Observation  
Benchmarking

## ABSTRACT

Remote sensing composed image retrieval (RSCIR) enables search in large satellite image archives using composed queries that combine a reference image with a textual modifier. Although RSCIR offers a flexible interface for expressing targeted retrieval intent, the transferability of modern composition methods to Earth observation (EO) imagery and their relevance to operational EO workflows remain underexplored. We address this gap through a unified benchmark and an application-oriented study. First, we systematically adapt and evaluate representative composed image retrieval methods with six vision-language backbones on PatternCom under a standardized protocol, analyzing their behavior across backbones, composition strategies, and query types. Second, we introduce xView2-CIR, a change-centric dataset for disaster and damage monitoring, where retrieval is conditioned on scene identity and a target post-event state. Our results show that training-free composition methods provide strong and scalable baselines for EO retrieval, while change-centric retrieval presents different challenges from attribute-based retrieval, particularly due to the need to preserve scene identity. Overall, this study establishes a practical benchmark for RSCIR and positions composed retrieval as a complementary tool for remote sensing image retrieval, archive exploration, and change analysis.

## 1. Introduction

The explosive growth of earth observation (EO) data has enabled large-scale monitoring of the planet, but has also made it increasingly difficult for users to navigate massive satellite image archives and retrieve content matching specific information needs. Remote sensing image retrieval (RSIR) [1] addresses this challenge by searching EO archives from a query, most commonly an image. Prior RSIR research spans unisource and cross-source [110] settings, as well as single-label [30, 80, 14] and multi-label [37, 34, 75] formulations. Despite this progress, a fundamental limitation persists: most RSIR systems are driven by a *single modality* query. A user must typically choose between an image query or a text query, which *restricts expressivity* when the retrieval intent includes *targeted modifications* beyond overall similarity.

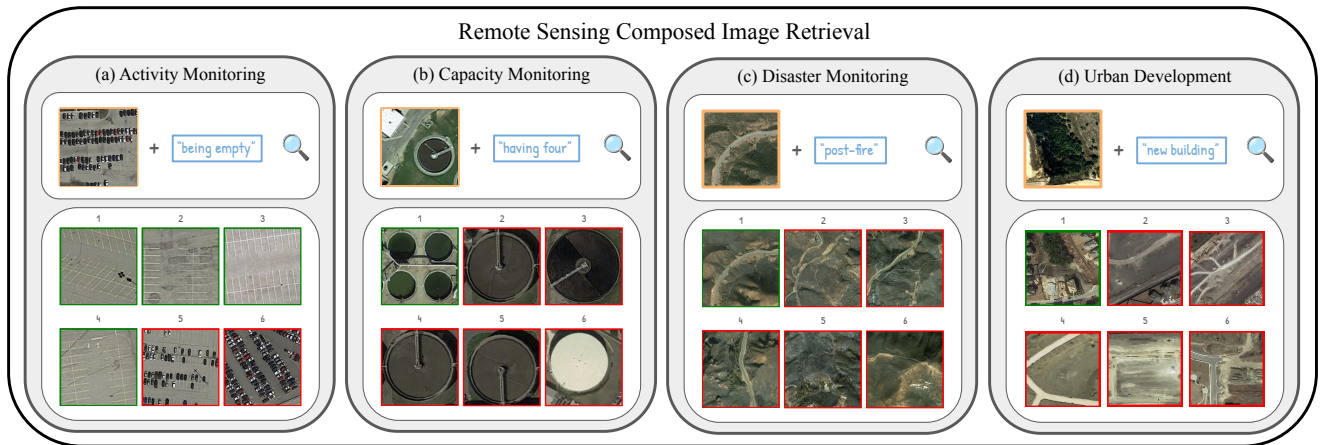
Many operational remote sensing workflows are inherently *compositional*. For example, an analyst may wish to retrieve locations visually *similar* to a reference region but with a different attribute, such as a parking lot with lower occupancy or a facility with a larger number of storage tanks, as shown in Figure 1(a,b). In change-centric applications, the goal may instead be to retrieve imagery of the *same* location under a different state, such as after a wildfire or recent construction, as shown in Figure 1(c,d). Such intent is difficult to express with unimodal queries alone. Remote sensing composed image retrieval (RSCIR) [64]

addresses this limitation by combining a reference image with a natural-language modifier, aiming to retrieve images that preserve the relevant visual context while satisfying the textual specification [71, 19, 65].

Despite its promise, RSCIR remains insufficiently studied from an applied EO perspective. Existing work [64] introduced the task, but it remains unclear how modern composed image retrieval methods behave under a unified evaluation protocol, how well they transfer across general-purpose and remote-sensing vision–language backbones, and whether composed retrieval remains effective when relevance depends on scene identity and a target post-event state rather than class-level attribute matching. These questions are important for assessing whether RSCIR can support practical EO workflows such as archive exploration, infrastructure monitoring, disaster response, and urban change analysis.

This work advances RSCIR through a unified benchmark and an application-oriented study. We first benchmark representative composition methods on PatternCom [64], an attribute-modification testbed derived from PatternNet [112]. We evaluate six vision–language backbones, including CLIP [67], SigLIP [105], RemoteCLIP [51], and SkyCLIP [94], and adapt representative composition methods such as Pic2Word [71], SEARLE [2], FreeDom [19], MagicLens [106], and BASIC [65] to EO imagery. Beyond aggregate retrieval metrics, we analyze method behavior across backbones, query types, and composition strategies to provide practical guidance for selecting RSCIR pipelines.

ORCID(s): 0000-0001-5381-0312 ( Bill Psomas)



**Figure 1: Composed Image Retrieval for applied Earth Observation.** We illustrate how composed queries (query image + query text) enable controllable retrieval by specifying a targeted change. (a) Activity monitoring: a parking-lot query image composed with being empty retrieves visually *similar* parking areas with low vehicle occupancy. (b) Capacity monitoring: A facility query image composed with having four retrieves visually *similar* sites exhibiting a different quantity attribute. (c) Disaster monitoring: a pre-event query image composed with post-fire retrieves post-event imagery of the *same* scene consistent with wildfire impact. (d) Urban development: a pre-construction query image composed with new buildings retrieves the *same* scene exhibiting recent construction. For each query, retrieved candidates (ranked left-to-right, top-to-bottom) are produced by FreeDOM [19] with OpenAI CLIP [67]; green borders indicate correct retrievals under the respective evaluation criterion (a,b: same class + target attribute value, c,d: same scene + target change), while red borders indicate mismatches. Images in (a,b) are from PatternNet [112], (c) from xView2 [41], and (d) from LEVIR-CC [50].

57 To connect benchmarking with operational EO needs, 58 we further introduce xView2-CIR, a new dataset derived 59 from xView2 [41] for change-centric composed retrieval in 60 disaster and damage monitoring. In xView2-CIR, a query 61 consists of a pre-event reference image and a textual modifier 62 such as post-fire, and the goal is to retrieve the corre- 63 sponding post-event image of the same location. This setting 64 evaluates whether composed retrieval can support location- 65 aware and semantically steerable search, complementing 66 established change-centric workflows such as change detec- 67 tion, damage assessment, and rapid mapping.

68 In summary, we make the following contributions:

- 69 1. We establish a unified benchmark for RSCIR on Pat- 70 ternCom, with domain-grounded adaptations of rep- 71 resentative composition methods and a standardized 72 protocol spanning six vision-language backbones.
- 73 2. We introduce xView2-CIR, a new evaluation dataset 74 for change-centric composed retrieval in disaster and 75 damage monitoring, where relevance depends on both 76 scene identity and target post-event state.
- 77 3. We provide empirical analysis and practical guidance 78 showing when composed retrieval is useful for remote 79 sensing applications, and how its behavior differs be- 80 tween attribute-based and change-centric settings.

## 81 2. Related Work

82 *Remote Sensing Image Retrieval.* With the aim to ef- 83 fectively *search* and *retrieve* information from extensive 84 remote sensing (RS) image archives, remote sensing image

85 retrieval (RSIR) can be categorized into *unisource* and *cross-* 86 *source* [110]. Initially, RSIR methods focus on handcrafted 87 and low-level visual features [58, 56, 63, 45, 5, 88, 74, 89, 9, 17]. With the advent of deep learning, neural networks 88 are utilized for unisource *single-label* retrieval: (a) as feature 89 extractors [47, 31, 6, 100, 59, 22, 83, 33, 70, 30], (b) 90 trained from scratch [111, 107, 113, 54, 92, 80, 93, 10], (c) 91 integrating attention modules [90, 91, 96, 11] and (d) using 92 metric learning [109, 8, 14, 53, 20, 52]. Neural networks are 93 also used for unisource *multi-label* [9, 37, 79, 81, 15, 34, 75], 94 cross-source *cross-sensors* [46, 57, 98, 97], cross-source 95 *cross-modal* [11, 99, 82, 78, 55, 103, 102] and cross-source 96 *cross-view* retrieval [32, 104, 85, 48, 39, 76]. Our work fills 97 a notable gap and enhances user intent expression in RSIR 98 by combining image with text.

99 *Composed Image Retrieval.* Image-to-image [66, 23, 61] 100 and text-to-image [72, 108, 21] retrieval provide ways to 101 explore large image archives. However, the most accurate 102 and flexible way to express the user intent is a query *com-* 103 *posed* of both an image and a text. Composed image retrieval 104 (CIR) [87, 13, 29, 3, 42, 71] aims to retrieve images not 105 only visually similar to the query image, but also altered to 106 align with the specifics of the query text. Traditionally, CIR 107 methods are supervised by *triplets* of the form *query image,* 108 *query text, target image* [87, 12, 101, 29, 68, 13, 42, 18]. The 109 labor-intensive process of labeling such triplets limits early 110 works to specific applications in fashion [28, 4, 95], physical 111 states [35], object attributes and composition [87, 49, 60]. 112 The emergence of vision-language models (VLMs) such as 113 CLIP [67], ALIGN [36], or BLIP [44] has recently enabled 114

115 *zero-shot composed image retrieval (ZS-CIR)*: one can com- 168  
 116 pose a query by manipulating embeddings alone, without 169  
 117 any task-specific training. Early ZS-CIR attempts include 170  
 118 Pic2Word [71] and SEARLE [2], which invert the visual em- 171  
 119 bedding back to text by either pre-training a small decoder 172  
 120 or optimizing a pseudo token at test-time respectively. More 173  
 121 recent work pushes this idea further with memory-based 174  
 122 inversion (FreeDom [19]) or caption-and-LLM pipelines 175  
 123 (CIReVL [38]). CompoDiff [25] casts the composition as 176  
 124 a diffusion [69] problem, gradually blending visual and 177  
 125 textual semantics. MagicLens [106] revisits triplet super- 178  
 126 vision at web scale, where a VLM with extra attention 179  
 127 layers is fine-tuned to project the image-text input to a 180  
 128 single embedding. Despite this rapid progress, the extent 181  
 129 to which ZS-CIR methods transfer to EO imagery and how 182  
 130 they should be benchmarked and used in real-world remote 183  
 131 sensing workflows remains insufficiently understood. This 184  
 132 paper addresses this gap through a unified benchmark, com- 185  
 133 plemented by an application-oriented study.

134 *Vision-Language Models*. Foundational VLMs such 187  
 135 as CLIP [67], ALIGN [36] and BLIP/BLIP-2 [44, 43] 188  
 136 are pre-trained on hundreds of millions, or even billions, 189  
 137 of web image-text pairs (e.g. LAION-2B [73]). Their 190  
 138 joint embedding spaces enable strong zero-shot transfer 191  
 139 to a wide range of tasks, including open-vocabulary 192  
 140 classification [67], detection [26], segmentation [77], and 193  
 141 captioning [84]. OpenCLIP [16] re-implements CLIP 194  
 142 and releases checkpoints trained directly on LAION-2B. 195  
 143 SigLIP [105] replaces the softmax contrastive loss with 196  
 144 a sigmoid variant and is trained on the multilingual 197  
 145 WebLI corpus. Several works adapt CLIP to the remote 198  
 146 sensing domain: RemoteCLIP [51] fine-tunes OpenAI 199  
 147 CLIP on caption-augmented RS datasets, CLIP LAION-RS 200  
 148 uses a 726k RS subset of LAION-2B (LAION-RS), and 201  
 149 SkyCLIP [94] leverages SkyScript, a million-scale image- 202  
 150 text pair dataset, constructed by linking satellite images 203  
 151 from Google Earth Engine [24] with OpenStreetMap [27] 204  
 152 annotations. In this work, we benchmark various VLMs 205  
 153 under identical composed image retrieval protocols, 206  
 154 revealing how generic versus RS-specialised pre-training 207  
 155 affects the performance on the proposed RSCIR task.

### 156 3. Task & Methodology

#### 157 3.1. Problem formulation

158 In remote sensing composed image retrieval (RSCIR), 212  
 159 the query consists of a reference image  $y$  and a textual 213  
 160 modifier  $t$ , denoted as  $q = (y, t)$ . The goal is to rank images 214  
 161  $x \in X$  from a database according to a composed similarity 215  
 162 score  $s(q, x) \in \mathbb{R}$ . Depending on the application, the query 216  
 163 image  $y$  is associated with either a semantic *class* label  $C_y$  217  
 164 or a *scene* label  $S_y$  identifying a geographic location. We 218  
 165 denote by  $A_y$  the attribute or state visible in  $y$ , and by  $A_t$  the 219  
 166 target attribute or state specified by the text.

167 We consider two relevance protocols:

- **Same class + target attribute.** A result  $x$  is relevant if it depicts the same class as the query and matches the target attribute, i.e.,  $C_x = C_y$  and  $A_x = A_t$ . This setting captures attribute-oriented search, such as retrieving similar facilities with a different color, shape, density, or quantity.
- **Same scene + target state.** A result  $x$  is relevant if it corresponds to the same scene/location as the query and matches the target state or change specified by the text, i.e.,  $S_x = S_y$  and  $A_x = A_t$ . This setting captures location-aware change-centric retrieval, such as retrieving the post-event image of the same site after a wildfire or flood.

The two protocols reflect different *operational goals*. The first evaluates class-conditioned attribute modification, while the second evaluates identity-preserving state retrieval. The latter is not a minor variant of the former: preserving scene identity while satisfying a target state may favor different composition mechanisms.

To define  $s$ , we make use of pre-trained VLMs that consist of a *visual encoder*  $f : \mathcal{I} \rightarrow \mathbb{R}^d$  and a *text encoder*  $g : \mathcal{T} \rightarrow \mathbb{R}^d$ , which map input images from image space  $\mathcal{I}$  and words from the text space  $\mathcal{T}$  to the same embedding space with dimension  $d$ . We extract the visual embedding  $\mathbf{v}_y = f(y) \in \mathbb{R}^d$  and the text embedding  $\mathbf{v}_t = g(t) \in \mathbb{R}^d$  to use as queries. Finally, the embedding of a dataset image  $x \in X$  is denoted as  $\mathbf{v}_x = f(x) \in \mathbb{R}^d$ . All embeddings are  $\ell_2$ -normalized.

#### 195 3.2. Baselines

196 *Unimodal*. These baselines score each database image us- 197  
 198 ing only one component of the composed query. Using the 199  
 199 VLM encoders  $f$  and  $g$ , we define the *text-only* baseline 200  
 200 as  $s_g(q, x) = g(t)^\top f(x)$  and the *image-only* baseline as 201  
 201  $s_f(q, x) = f(y)^\top f(x)$ . These baselines isolate the contribu- 202  
 202 tion of the textual modifier and visual reference, respectively.

203 *Multimodal*. These baselines combine the two unimodal 204  
 204 scores by fusing the image- and text-based similarities. A 205  
 205 simple fusion is *averaging*:  $s_a(q, x) = \frac{1}{2}[s_g(q, x) + s_f(q, x)]$ , 206  
 206 which yields the same ranking as *summing* the similarities. 207  
 207 We denote this baseline as *text + image*. While straight- 208  
 208 forward, this fusion can be biased toward the image signal 209  
 209 in practice, since same-modality similarities (image-image) 210  
 210 often have a different scale than cross-modal similarities 211  
 211 (text-image), even under  $\ell_2$  normalization. A second fusion 212  
 212 computes similarity by *multiplying* the unimodal scores as 213  
 213  $s_m(q, x) = s_f(q, x) \times s_g(q, x)$ . This method emphasizes 214  
 214 retrieval results with high agreement between modalities, 215  
 215 inherently penalizing disagreement, and serving as a form of 216  
 216 *soft normalization* that mitigates bias toward either modality. 217  
 217 We denote this baseline as: *text  $\times$  image*.

#### 217 3.3. Methods

218 Beyond unimodal and score-fusion baselines, we adapt 219  
 219 and evaluate representative composed image retrieval meth- 220  
 220 ods from computer vision. These methods cover textual

221 inversion, memory-based query construction, caption-and-  
 222 LLM composition, diffusion-based embedding composition,  
 223 supervised multimodal pooling, and training-free feature  
 224 calibration.

225 *WeiCom* [64] is a training-free score-fusion method. It  
 226 computes image-to-image and text-to-image similarities,  
 227 standardizes each score distribution over the database, maps  
 228 the standardized scores to  $[0, 1]$  using the Gaussian CDF, and  
 229 combines them as  $s_{WC}(q, x) = \lambda s'_g(q, x) + (1 - \lambda)s'_f(q, x)$ ,  
 230 where  $s'_f$  and  $s'_g$  are the calibrated image and text scores,  
 231 and  $\lambda \in [0, 1]$  controls the modality weight.

232 *Pic2Word* [71] casts composition as *textual inversion*: it  
 233 learns an approximate inverse of the text encoder, mapping  
 234 an image embedding back into a *text-space token representa-*  
 235 *tion* that can be composed with the modifier. Concretely,  
 236 given  $\mathbf{v}_y = f(y)$ , it predicts a pseudo-text representation  
 237  $y^* \approx g^{-1}(\mathbf{v}_y)$  with a small network trained on large-scale  
 238 image-text pairs. The composed query is then formed in text  
 239 space via concatenation and encoded as  $g([y^*; t])$ , which is  
 240 used for text-to-image retrieval.

241 *SEARLE* [2] also follows textual inversion, but avoids  
 242 pretraining by performing *test-time optimization*. For each  
 243 query image, it directly optimizes a pseudo token  $y^*$  such  
 244 that  $g(y^*)$  matches  $\mathbf{v}_y$  under a cosine objective, comple-  
 245 mented by an LLM/GPT-based [7] regularizer that encour-  
 246 ages linguistic plausibility. The optimized  $y^*$  is then con-  
 247 catenated with the modifier and encoded with  $g(\cdot)$  for re-  
 248 trieval. This per-query optimization is computationally heav-  
 249 ier, but can yield strong inversion quality.

250 *FreeDom* [19] replaces explicit inversion with a *mem-*  
 251 *ory* of textual anchors. It pre-encodes a vocabulary  $W =$   
 252  $\{w_i\}_{i=1}^N$  through the text encoder  $g(\cdot)$  to obtain  $V_W =$   
 253  $\{g(w_i)\}_{i=1}^N$  and performs nearest-neighbor search from the  
 254 image embedding  $\mathbf{v}_y$  to retrieve a set of words whose em-  
 255 beddings best match  $\mathbf{v}_y$ . Since the mapping between words  
 256 and embeddings is explicit, this provides an interpretable  
 257 textual surrogate for the image query without optimization  
 258 or pretraining. The retrieved words are concatenated with the  
 259 modifier and encoded with  $g(\cdot)$  for text-to-image retrieval.  
 260 *FreeDom* further improves robustness through query expan-  
 261 sion using visually similar images.

262 *CIReVL* [38] follows a *caption-and-LLM* pipeline. It con-  
 263 verts the image query to natural language via captioning  
 264 rather than inversion. A captioner (e.g., BLIP/BLIP-2 [44,  
 265 43]) produces an image description, which is then merged  
 266 with the modifier using an LLM to form a refined query sen-  
 267 tence. Retrieval is performed as text-to-image matching with  
 268 the VLM. This pipeline favors human-readable intermediate  
 269 representations and can benefit from strong captioning and  
 270 rewriting quality.

271 *CompoDiff* [25] recasts composition as a *diffusion* process  
 272 in the joint VLM embedding space. Starting from the image

273 embedding  $\mathbf{v}_y$ , a denoising diffusion model gradually injects  
 274 the textual condition  $t$  (with image self-conditioning), pro-  
 275 ducing a composed embedding  $\tilde{\mathbf{v}}_{y,t}$  after  $T$  steps. Retrieval  
 276 is then performed with cosine similarity  $\tilde{\mathbf{v}}_{y,t}^\top \mathbf{v}_x$ .

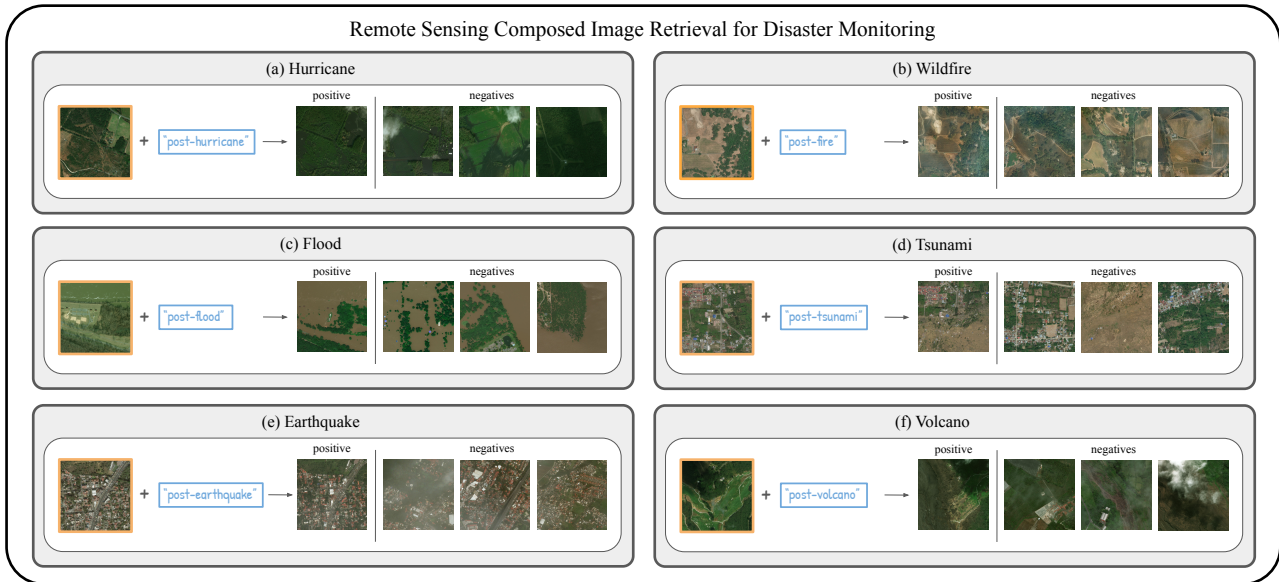
277 *MagicLens* [106] follows a *supervised multimodal pool-*  
 278 *ing* pipeline. It learns a dedicated composition module via  
 279 supervision. It is trained on web-mined triplets (query im-  
 280 age, query text, target image), and fine-tunes a VLM aug-  
 281 mented with additional attention/pooling layers to map the  
 282 joint image-text input to a single embedding. Database  
 283 images are mapped to the same space and matched by cosine  
 284 similarity.

285 *BASIC* [65] is a *training-free compositional scoring*  
 286 method operating on frozen VLM features. It starts from  
 287 the standard image-to-image and text-to-image similarities,  
 288 and aims to make their combination more reliable by  
 289 (i) reducing modality- and domain-specific bias, (ii)  
 290 suppressing nuisance directions in the embedding space, and  
 291 (iii) fusing the two modalities with an “AND”-like criterion.  
 292 Concretely, *BASIC* first *centers* image and text embeddings  
 293 by subtracting mean features (optionally estimated from  
 294 a large external corpus), which mitigates global bias in  
 295 the representation. It then applies a *semantic projection*  
 296 derived from two textual corpora: a positive/object corpus  
 297  $C_+$  and a negative/style/context corpus  $C_-$ . A contrastive  
 298 PCA-style construction emphasizes directions that are  
 299 informative for object content while de-emphasizing  
 300 directions associated with style, viewpoint, or acquisition  
 301 context. Optionally, *BASIC* further strengthens the visual  
 302 cue via *query expansion* using top retrieved neighbors.  
 303 Finally, *BASIC* computes modality scores against the  
 304 database and applies a lightweight *score calibration*  
 305 to reduce scale mismatch between modalities, before  
 306 combining them with a multiplicative fusion regularized  
 307 by a Harris-like penalty.

## 308 4. Experiments

### 309 4.1. Experimental setup

310 *Datasets*. We conduct our study on three datasets that  
 311 cover both *attribute-driven* and *change-centric* composed  
 312 retrieval in EO imagery. *PatternCom* [64] is an existing  
 313 RSCIR benchmark derived from *PatternNet* [112], a high  
 314 resolution dataset with 38 classes (800 images per class,  
 315 256×256 pixels). *PatternCom* organizes *PatternNet* into  
 316 composed queries of the form (query image, query text),  
 317 where the text specifies a target *attribute value* for the  
 318 class depicted in the query image (e.g., shape: oval for  
 319 swimming pool). It spans six attribute types (*color, context,*  
 320 *density, existence, quantity, shape*) with class-specific value  
 321 sets and highly variable numbers of positives per query.  
 322 In *PatternCom*, relevance is defined by the *same class*  
 323 + *target attribute* criterion. As illustrated in [Figure 1](#)  
 324 (a,b), this protocol naturally supports attribute-oriented



**Figure 2: Remote sensing composed image retrieval for disaster monitoring.** We restructure xView2 [41] into a composed retrieval setting where a query consists of a *pre-event* reference image of a specific location and a *textual modifier* describing the desired post-event state (e.g., post-hurricane). For each disaster type (a–f), we illustrate the query and the corresponding relevance criterion: *positives* are images of the *same scene/location* depicting the specified post-event state, while *negatives* are non-matching candidates (e.g., different locations with the target change or the same location without the target change).

325 RS use cases such as activity/occupancy monitoring and  
 326 capacity/infrastructure monitoring.

327 PatternCom does not fully reflect operational scenarios  
 328 where users seek *the same location* under a different state.  
 329 To bridge this gap, we introduce xView2-CIR, a new dataset  
 330 built by restructuring xView2 [41] into a composed retrieval  
 331 setting for disaster and damage monitoring. In xView2-  
 332 CIR, each query pairs a *pre-event* reference image of a geo-  
 333 registered scene with a text modifier describing the desired  
 334 *post-event* state (e.g., post-hurricane, post-fire); relevance  
 335 follows the *same scene/location + target state/change* crite-  
 336 rion, as depicted in Figure 2. The unique positive of each  
 337 query is the post-event image of the same geo-registered  
 338 location. All remaining images are treated as negatives, in-  
 339 cluding images from different locations that exhibit the same  
 340 disaster category and images from the same location that do  
 341 not match the target post-event state, when applicable. This  
 342 construction isolates the challenge of identity-preserving  
 343 change retrieval.

344 Finally, to further highlight the applicability of RSCIR to  
 345 urban change analysis, we include qualitative examples from  
 346 LEVIR-CC [50], which provides paired imagery for building  
 347 change detection. We use a small, partially structured subset  
 348 to form composed queries for urban development (e.g.,  
 349 new building), following the *same scene/location + target*  
 350 *change* criterion; these examples are used for qualitative  
 351 analysis only (e.g., Figure 1 (d)), and are not included in the  
 352 quantitative benchmark due to the absence of a fully curated  
 353 CIR protocol for this dataset in our study. A comprehensive  
 354 overview of the dataset statistics is provided in A.1.

355 *Networks.* We evaluate six vision-language models,  
 356 all using the ViT-L/14 vision backbone: CLIP [16]  
 357 LAION-2B [73], RemoteCLIP [51], OpenAI CLIP [67],  
 358 SigLIP [105], CLIP LAION-RS [94], and SkyCLIP-50 [94].  
 359 We provide further details on the networks in A.2.

360 *Vocabularies.* Recent CIR methods [2, 19, 65] leverage  
 361 a vocabulary of possible categories or concepts. However,  
 362 many of these, such as Open Images v7 [40] with 21k  
 363 concepts, originate from general-purpose computer vision  
 364 datasets. To better align these methods with remote sensing  
 365 semantics, we create a family of synthetic vocabularies using  
 366 GPT-4o [62], which we will distribute to avoid redistribution  
 367 of any copyrighted content. These vocabularies, denoted as  
 368 RSText, consist of increasing sizes: 150 to two thousand dis-  
 369 tinct concepts. The generation process (presented in subsec-  
 370 tion A.3) follows a structured prompt that emphasizes fine-  
 371 grained, diverse, and remote-sensing-relevant terminology  
 372 across a wide range of thematic areas. To combine both  
 373 general-purpose and domain-specific semantics, we further  
 374 merge RS-Text-2k with the Open Images v7, resulting in a  
 375 hybrid vocabulary of 23k classes, referred to as HybridText-  
 376 23k.

377 *Evaluation protocols.* For a query  $q$ , we compute Av-  
 378 erage Precision (AP) as the mean of the precision values  
 379 at the ranks where relevant items are retrieved. We then  
 380 aggregate AP across queries to obtain mAP, which reflects  
 381 both relevance and ranking quality. For PatternCom, we  
 382 follow an *attribute-balanced* protocol: AP is first averaged  
 383 over queries within each attribute type (e.g., *color*), and the  
 384 final score is the unweighted average across attribute types.

385 This corresponds to a *macro-averaged mAP* ( $mAP_{macro}$ ),  
 386 ensuring that each attribute contributes equally. For xView2-  
 387 CIR, the query distribution across disaster types is im-  
 388 balanced. Therefore, we report both (i) the same *macro-*  
 389 *averaged mAP* over disaster categories,  $mAP_{macro}$ , which  
 390 weights each disaster type equally, and (ii) the standard  
 391 *overall mAP*,  $mAP_{overall}$ , computed by averaging AP over  
 392 *all queries* irrespective of category. Reporting both provides  
 393 a fair view of performance:  $mAP_{macro}$  reflects robustness  
 394 across rare disaster types, while  $mAP_{overall}$  reflects expected  
 395 performance under the natural query frequency.

396 *Hyperparameters.* Both PatternCom and xView2-CIR are  
 397 used strictly as evaluation sets. To avoid tuning on test  
 398 data, we keep method hyperparameters fixed to the default  
 399 settings reported in the original papers (or their official  
 400 implementations) whenever applicable. When a method ex-  
 401 poses several knobs (*e.g.*, visual/textual expansion sizes,  
 402 inversion iterations, diffusion steps), we do *not* optimize  
 403 them on PatternCom or xView2-CIR. Instead, we report a  
 404 sensitivity analysis in the ablation section to characterize  
 405 how performance varies with key hyperparameters, and we  
 406 use a single fixed configuration for all main comparisons. For  
 407 instance, for WeiCom we report the  $\lambda$  sensitivity separately,  
 408 while using a single fixed  $\lambda$  in the main tables for fair  
 409 comparison.

410 **4.2. Experimental analysis**

411 *Benchmark.* Table 1, 2, and 3 report the main quantita-  
 412 tive results on PatternCom and xView2-CIR, respectively,  
 413 covering a broad range of composition methods and vision-  
 414 language models. We include unimodal baselines (*text-only*,  
 415 *image-only*) and two simple multimodal fusions: score av-  
 416 eraging (*text + image*) and score multiplication (*text ×*  
 417 *image*). We further benchmark representative state-of-the-  
 418 art CIR methods from the computer vision literature: Com-  
 419 poDiff [25], Pic2Word [71], SEARLE [2], CIREVL [38],  
 420 MagicLens [106], WeiCom [64], BASIC [65], and Free-  
 421 Dom [19], after applying domain-grounded adaptations re-  
 422 quired for EO imagery. Our adaptation principle is intention-  
 423 ally conservative: for each method, we preserve the original  
 424 inference mechanism and introduce only minimal domain-  
 425 grounded modifications needed to make it meaningful for  
 426 EO imagery, such as replacing generic vocabularies, ad-  
 427 justing prompt templates, or using remote-sensing-relevant  
 428 textual resources. We generally avoid benchmark-specific  
 429 tuning beyond what is required for compatibility. Detailed  
 430 adaptation choices and implementation details are provided  
 431 in A.4.

432 **4.3. Experimental results**

433 *Quantitative results on PatternCom* Table 1 and 2 sum-  
 434 marize attribute-modification retrieval performance on Pat-  
 435 ternCom across six VLM backbones. Unimodal baselines  
 436 are consistently weak: *image-only* usually outperforms *text-*  
 437 *only*, confirming that visual similarity alone is useful but  
 438 insufficient for satisfying the textual modifier. Simple multi-  
 439 modal fusion provides clear gains over unimodal retrieval. In

**Table 1**  
**Attribute modification performance on PatternCom (mAP, %).** Results for three vision—language backbones and up to eight composition methods. We report *macro-averaged mAP* ( $mAP_{macro}$ ): for each attribute type, AP is averaged over all queries targeting a given attribute value (*e.g.*, rectangular for Shape), then averaged over the remaining values of the same attribute type (*e.g.*, oval, kidney-shaped), and finally averaged across attribute types. Avg.: resulting  $mAP_{macro}$  over all attribute types and classes; **bold**: best; underline: second.

(a) CLIP LAION-2B							
METHOD	Color	Context	Density	Existence	Quantity	Shape	Avg.
Text-only	13.47	4.83	3.58	2.00	3.31	6.22	5.57
Image-only	14.66	8.32	13.49	16.47	7.84	15.76	12.74
Text + Image	23.13	11.02	15.87	16.93	10.13	21.38	16.41
Text × Image	40.97	11.87	14.48	14.36	20.31	23.99	21.00
SEARLE	14.75	7.98	13.63	15.23	8.01	15.86	12.58
CIREVL	17.79	<u>28.55</u>	<u>17.69</u>	<u>35.14</u>	14.95	<u>25.91</u>	<u>23.34</u>
WeiCom	46.08	17.45	16.49	8.36	<u>18.15</u>	23.97	21.75
BASIC	29.26	6.38	15.29	18.54	12.18	17.18	16.47
FreeDom	<b>46.55</b>	<b>43.49</b>	<b>21.32</b>	<b>46.98</b>	<b>25.09</b>	<b>48.13</b>	<b>38.59</b>

(b) RemoteCLIP							
METHOD	Color	Context	Density	Existence	Quantity	Shape	Avg.
Text-only	10.75	8.87	22.16	6.98	8.25	24.12	13.52
Image-only	14.40	6.62	15.11	13.10	6.99	15.18	11.90
Text + Image	23.67	10.01	18.45	13.98	7.97	19.63	15.62
Text × Image	47.20	19.65	27.09	12.97	14.59	40.60	27.02
CompoDiff	8.58	17.88	14.41	6.04	9.34	12.46	11.45
Pic2Word	40.88	<b>40.26</b>	17.41	16.92	9.18	27.98	25.44
SEARLE	14.44	6.00	13.49	12.95	7.29	14.86	11.51
CIREVL	42.80	36.79	21.34	<b>43.81</b>	19.35	<u>35.55</u>	<u>33.27</u>
WeiCom	43.68	31.45	<b>39.94</b>	14.92	20.51	29.78	30.05
BASIC	47.65	15.45	17.55	22.03	20.76	25.00	24.74
FreeDom	<b>49.80</b>	<b>38.80</b>	<b>28.31</b>	36.64	<b>26.77</b>	<b>39.49</b>	<b>36.64</b>

(c) OpenAI CLIP							
METHOD	Color	Context	Density	Existence	Quantity	Shape	Avg.
Text-only	3.49	7.75	2.51	3.31	1.01	2.94	3.50
Image-only	13.61	7.86	18.19	13.30	8.30	15.08	12.72
Text + Image	16.59	11.13	19.88	15.13	8.96	17.55	14.87
Text × Image	27.59	14.65	15.63	17.97	9.08	17.77	17.12
CompoDiff	10.70	8.32	16.83	12.36	4.67	13.50	11.06
Pic2Word	27.09	20.93	19.04	9.81	8.05	19.24	17.36
SEARLE	15.24	7.25	16.20	13.30	9.00	15.78	12.80
CIREVL	29.32	<u>31.14</u>	16.02	<u>37.23</u>	11.70	22.88	24.72
MagicLens	34.63	18.98	14.54	21.46	13.21	17.70	20.09
WeiCom	27.42	25.15	13.63	21.98	7.01	15.33	18.42
BASIC	38.31	24.53	26.35	30.94	<b>18.03</b>	<b>35.73</b>	28.98
FreeDom	<b>39.08</b>	<b>45.87</b>	<b>36.72</b>	<b>37.25</b>	<b>13.22</b>	<b>28.69</b>	<b>33.47</b>

440 particular, *text × image* is often stronger than *text + image*,  
 441 although it is less stable on some backbones, such as SigLIP.

442 Among composition methods, FreeDom is the strongest  
 443 performer, achieving the best average mAP across all six  
 444 backbones, with particularly high absolute performance  
 445 on SigLIP (46.61%) and strong results on RemoteCLIP,  
 446 SkyCLIP-50, and CLIP LAION-RS. CIREVL and BASIC  
 447 form the next tier. CIREVL performs particularly well on  
 448 remote-sensing-adapted CLIP backbones like RemoteCLIP  
 449 and SkyCLIP-50. BASIC is the second strongest with  
 450 OpenAI CLIP and SigLIP. WeiCom remains a competitive  
 451 lightweight baseline, while MagicLens shows that super-  
 452 vised multimodal heads trained on natural-image triplets

**Table 2**

**Attribute modification performance on PatternCom (mAP, %).** Results for three vision–language backbones and up to seven composition methods. We report *macro-averaged mAP* ( $mAP_{macro}$ ): for each attribute type, AP is averaged over all queries targeting a given attribute value (e.g., rectangular for Shape), then averaged over the remaining values of the same attribute type (e.g., oval, kidney-shaped), and finally averaged across attribute types. Avg.: resulting  $mAP_{macro}$  over all attribute types and classes; **bold**: best; underline: second.

(a) SigLIP							
METHOD	Color	Context	Density	Existence	Quantity	Shape	Avg.
Text-only	7.99	3.93	2.30	9.93	0.58	31.16	9.32
Image-only	15.03	10.67	21.46	21.13	8.94	17.37	15.77
Text + Image	25.11	11.46	20.47	21.31	10.66	27.92	19.49
Text × Image	18.80	5.96	6.57	18.77	1.28	43.66	15.84
WeiCom	24.77	<u>15.32</u>	9.06	22.37	4.59	43.64	19.96
BASIC	<u>36.11</u>	<u>10.30</u>	<u>21.76</u>	<u>24.68</u>	<u>20.43</u>	19.46	<u>22.12</u>
FreeDom	<b>47.99</b>	<b>49.88</b>	<b>22.84</b>	<b>54.70</b>	<b>37.05</b>	<b>67.17</b>	<b>46.61</b>

(b) CLIP LAION-RS							
METHOD	Color	Context	Density	Existence	Quantity	Shape	Avg.
Text-only	3.85	6.13	2.60	5.56	0.99	4.57	3.95
Image-only	13.35	7.51	15.82	15.79	7.50	15.91	12.65
Text + Image	20.94	10.07	16.54	16.91	8.75	20.78	15.67
Text × Image	21.38	12.98	10.91	16.91	10.20	22.32	15.78
CompoDiff	19.70	11.32	<u>17.98</u>	17.97	9.67	20.36	16.17
Pic2Word	31.18	19.58	16.53	11.87	7.88	19.03	17.68
SEARLE	14.89	10.23	14.34	17.20	8.31	16.51	13.58
CIReVL	32.57	36.45	15.33	34.68	<b>13.78</b>	24.88	<u>26.28</u>
WeiCom	29.63	30.45	11.19	21.13	8.60	<u>25.85</u>	21.14
BASIC	<u>35.51</u>	14.06	16.32	25.56	12.18	21.81	20.91
FreeDom	<b>41.78</b>	<b>51.34</b>	<b>21.88</b>	<b>38.22</b>	<u>13.29</u>	<b>33.58</b>	<b>33.35</b>

(c) SkyCLIP-50							
METHOD	Color	Context	Density	Existence	Quantity	Shape	Avg.
Text-only	4.60	8.44	3.44	5.61	1.02	7.34	5.08
Image-only	14.48	9.03	17.82	20.11	8.77	16.32	14.42
Text + Image	23.49	12.64	<u>19.54</u>	21.42	10.24	21.80	18.19
Text × Image	37.03	15.84	16.24	23.99	12.90	31.51	22.92
CompoDiff	20.59	12.67	17.57	21.98	10.23	22.31	17.56
Pic2Word	34.57	21.58	18.23	13.09	10.03	21.93	19.91
SEARLE	16.57	11.68	16.41	20.22	9.22	16.88	15.16
CIReVL	34.90	34.93	17.09	<u>35.03</u>	<u>14.39</u>	28.92	<u>27.54</u>
WeiCom	40.46	<u>38.10</u>	18.17	27.91	10.10	<u>31.52</u>	27.71
BASIC	34.89	12.20	16.64	28.98	12.32	18.30	20.56
FreeDom	<b>49.88</b>	<b>50.47</b>	<b>29.71</b>	<b>38.00</b>	<b>14.58</b>	<b>39.99</b>	<b>37.11</b>

can transfer to EO imagery to some extent. In contrast, continuous inversion and diffusion-based approaches are generally less competitive. Pic2Word and SEARLE yield modest gains at best, while CompoDiff is often unstable and can even underperform the *image-only* baseline. This suggests that methods relying on synthesized query embeddings or pseudo-token inversion are more sensitive to domain shift in EO imagery.

Overall, PatternCom favors methods that preserve class identity while injecting attribute-specific textual cues. Strong textual surrogates and domain-grounded vocabularies are therefore particularly beneficial, explaining the consistent advantage of FreeDom in this setting.

**Table 3**

**Disaster monitoring performance on xView2-CIR (mAP, %).** Results for two vision–language backbones and three composition methods. We report per-disaster mAP (%) for post- $\times$  modifiers (e.g., Hurricane). Avg.:  $mAP_{macro}$  across disasters (equal weight per disaster); Total:  $mAP_{overall}$  (weighted by the number of queries per disaster); **bold**: best; underline: second.

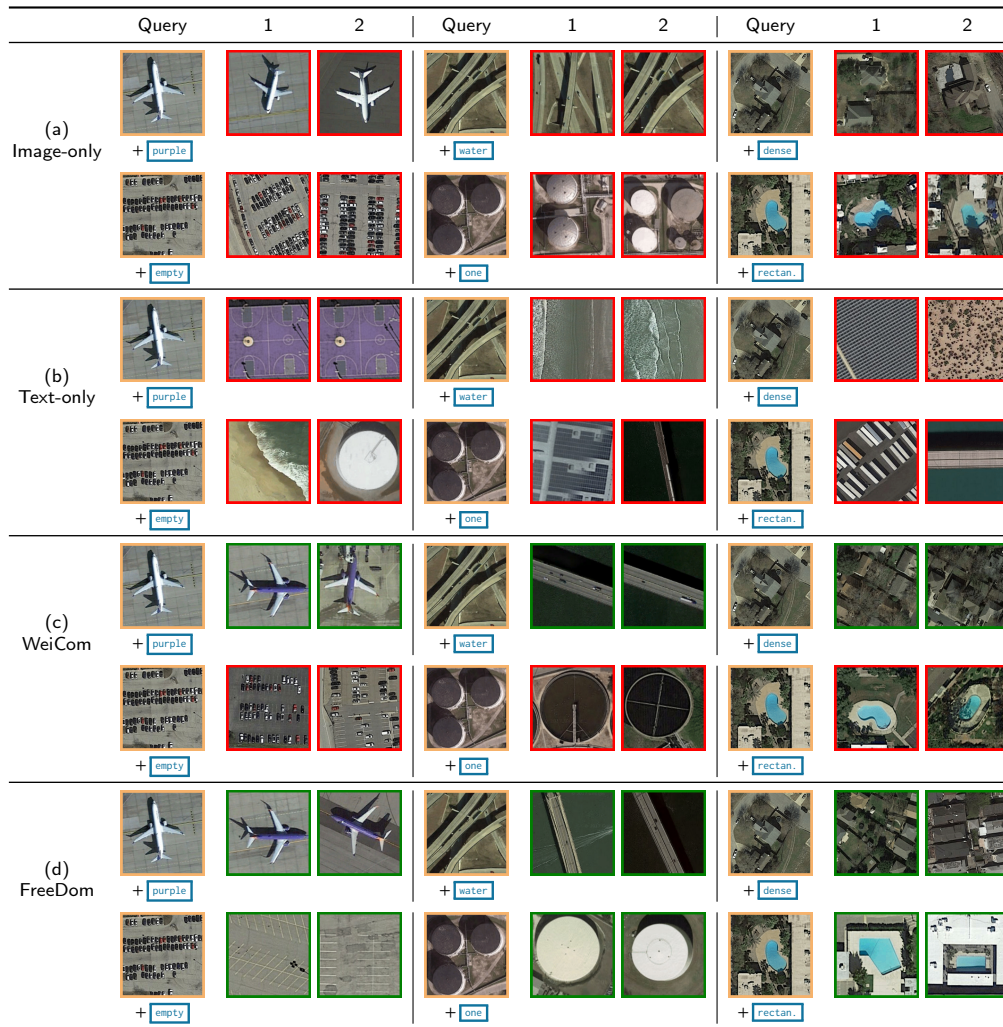
(a) OpenAI CLIP								
METHOD	Hurricane	Wildfire	Flood	Tsunami	Earthquake	Volcano	Avg.	Total
Text-only	2.50	1.61	3.25	18.08	4.19	15.77	7.57	2.97
Image-only	6.99	2.39	1.44	25.17	8.35	9.69	9.00	5.53
Text + Image	14.61	7.45	4.75	37.46	16.63	25.49	17.73	12.14
Text × Image	18.11	12.12	8.36	37.63	<u>17.93</u>	<b>61.46</b>	25.94	16.35
WeiCom	<b>25.94</b>	13.88	<u>9.43</u>	<u>42.96</u>	<b>36.48</b>	55.34	<b>30.67</b>	<b>21.40</b>
BASIC	22.71	<b>18.49</b>	<b>14.59</b>	<b>46.15</b>	1.69	59.20	27.14	21.38
FreeDom	18.01	<u>17.46</u>	3.77	35.03	3.21	25.88	17.23	16.84

(b) SigLIP								
METHOD	Hurricane	Wildfire	Flood	Tsunami	Earthquake	Volcano	Avg.	Total
Text-only	2.32	1.81	1.09	7.10	0.88	13.77	4.49	2.31
Image-only	10.65	9.74	10.68	24.15	14.44	20.77	15.07	10.97
Text + Image	15.86	19.94	13.05	26.29	<u>19.92</u>	26.17	20.20	17.50
Text × Image	10.74	13.69	<u>14.20</u>	17.54	<b>23.17</b>	53.96	22.22	13.08
WeiCom	<u>20.19</u>	<b>30.61</b>	9.95	<u>29.46</u>	18.56	32.09	<u>23.48</u>	<b>23.14</b>
BASIC	<b>21.54</b>	12.30	13.98	<b>36.52</b>	13.86	<b>56.79</b>	<b>25.83</b>	<u>18.52</u>
FreeDom	13.07	<u>15.37</u>	<b>21.22</b>	19.33	1.60	48.61	19.99	15.37

*Quantitative results on xView2-CIR* Table 3 reports performance on xView2-CIR under the *same scene/location + target state* criterion. Unimodal baselines remain weak, while score fusion provides clear gains across both backbones. *Text + image* is consistently reliable, whereas *text × image* is more variable across disaster types. Among composition methods, WeiCom is the most robust overall, achieving the best TOTAL score for both OpenAI CLIP and SigLIP, suggesting that calibrated modality fusion is well suited to identity-preserving change retrieval. In contrast, FreeDom and BASIC are less effective than on PatternCom, partly because their query-expansion mechanisms can introduce visually similar but geographically different scenes, which is harmful when relevance depends on scene identity. These results show that change-centric RSCIR poses different challenges from attribute-based retrieval and should be evaluated as a distinct setting.

*Qualitative results on PatternCom.* In Figure 3, we visualize representative retrieval outputs on PatternCom with OpenAI CLIP, comparing unimodal baselines (image-only, text-only) against two multimodal methods (WeiCom and FreeDom). The examples clearly illustrate the limitations of unimodal retrieval. WeiCom generally improves over unimodal baselines by leveraging both modalities, but it can still be sensitive to modality dominance depending on the attribute type. In the *activity monitoring* example (parking lot + empty), the retrieved results remain biased toward the visual layout of the reference image, yielding parking-lot-like structures that do not clearly satisfy the “empty” intent. In the *capacity monitoring* example (storage tanks + one), the textual cue can dominate, pulling results with fewer tanks that are only loosely aligned with the specific query instance. In contrast, FreeDom is consistently the most reliable across



**Figure 3: Qualitative composed retrieval results on PatternCom with OpenAI CLIP.** Comparison between unimodal and multimodal methods. Each query is shown as a **reference image** combined with a boxed **textual modifier**. Columns report the top-2 retrieved results. Retrieval is evaluated under the *same class + target attribute value* relevance criterion.

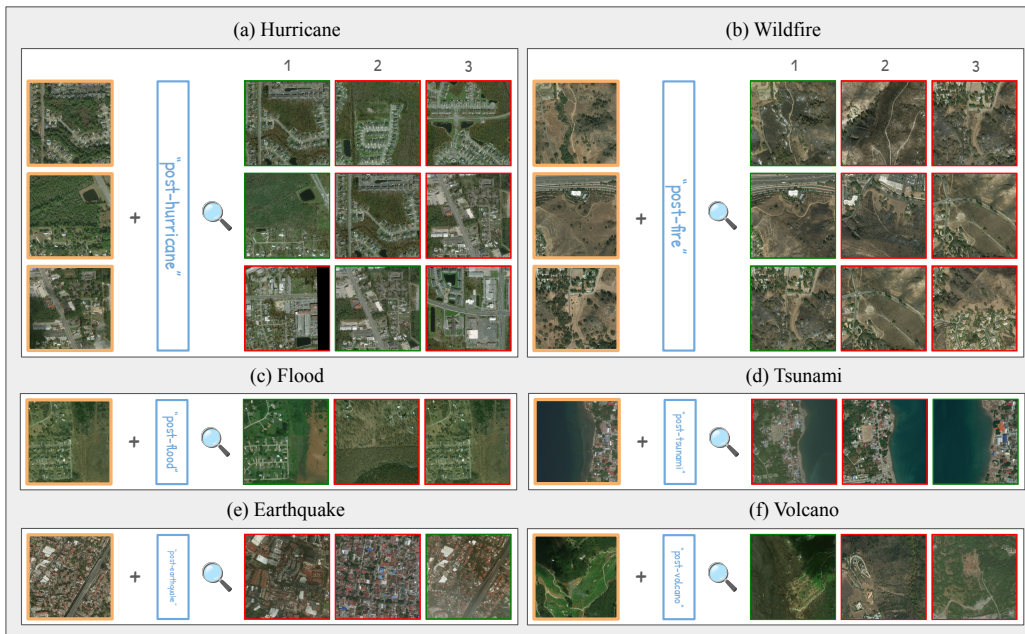
all shown cases, producing results that simultaneously respect the query class and the target attribute value, aligning well with PatternCom’s relevance criterion.

**Qualitative results on xView2-CIR.** In Figure 4, we visualize qualitative composed retrieval results on xView2-CIR using OpenAI CLIP with WeiCom. Unlike PatternCom, relevance in xView2-CIR requires *the same scene or location and the target state* (post-event), which exposes a key failure mode of generic composition: many retrieved images match the requested disaster state but drift to a different geographic area with similar visual cues. This is especially apparent for hurricanes and wildfires, where widespread debris/burn patterns and textured backgrounds can look plausible across locations, leading to several off-scene false positives.

Still, WeiCom succeeds in a subset of cases where the scene has distinctive geometry or landmarks (*e.g.*, road layouts, coastline structure, dense urban blocks), allowing the query image to anchor location while the text steers the retrieval toward the correct post-event state. Overall, the figure

highlights that xView2-CIR is less forgiving than attribute-only benchmarks: effective composed retrieval must jointly preserve *instance identity* (same place) while applying a *state change* constraint, and methods that lean too heavily on “state-like” appearance cues tend to retrieve the right disaster but the wrong scene.

**Qualitative results on LEVIR-CC.** To illustrate more applied EO use cases, Figure 5 shows qualitative results on LEVIR-CC [50] using WeiCom with OpenAI CLIP. As shown in the examples for *building construction* and *road construction*, the method takes a reference image together with a change-oriented modifier (*e.g.*, new building) and retrieves candidate scenes that match the requested change, under the *same location/scene + target state* relevance criterion. In these cases, the top retrieved results frequently correspond to plausible construction changes, suggesting that composed retrieval can serve as a lightweight interface for change-focused search in real remote sensing workflows (*e.g.*, rapid infrastructure monitoring).



**Figure 4: Qualitative composed retrieval results on xView2-CIR with OpenAI CLIP and WeiCom.** Each query combines a *pre-disaster* reference image with a boxed textual modifier (post-\*) indicating the target post-event state. We show the top retrieved *post-disaster* results per disaster type. Retrieval is evaluated under the *same scene/location + target state* relevance criterion, so visually plausible post-event matches from different locations are considered negatives.

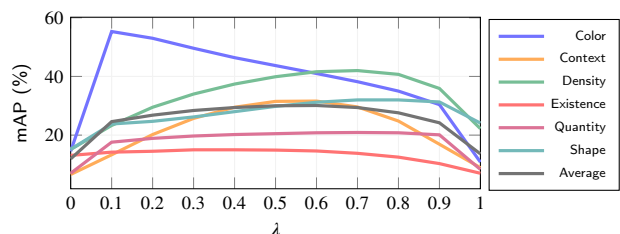


**Figure 5: Qualitative composed retrieval on LEVIR-CC with OpenAI CLIP and WeiCom.** Each query is shown as a reference image combined with a textual change modifier. We report the top retrieved candidates under the *same location + target state* relevance criterion.

#### 4.4. Sensitivity and Ablation Analysis

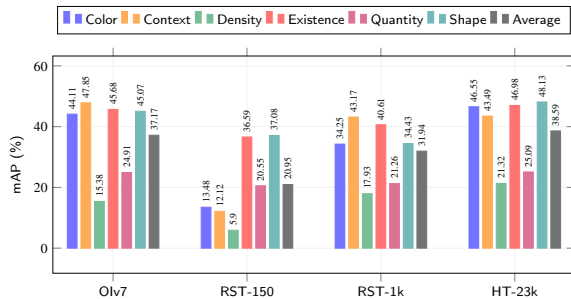
*Impact of  $\lambda$  in WeiCom.* Figure 6 analyzes the modality-control parameter  $\lambda$  of WeiCom with RemoteCLIP, where  $\lambda=0$  corresponds to *image-only* and  $\lambda=1$  to *text-only* retrieval. Although all main results use a fixed  $\lambda=0.5$  to avoid tuning on PatternCom, the sensitivity analysis shows that performance is maximized at intermediate values, with the best average mAP obtained at  $\lambda=0.6$ . Similar trends are observed across backbones, where the optimal  $\lambda$  consistently lies away from the unimodal extremes (e.g.,  $\lambda=0.3$  for CLIP LAION-2B). This broad mid-range optimum indicates that both the visual reference and textual modifier contribute

549 meaningfully to attribute-based retrieval. The main exception is *color*, which peaks at smaller  $\lambda$  values, suggesting  
550 stronger dependence on the visual signal.  
551



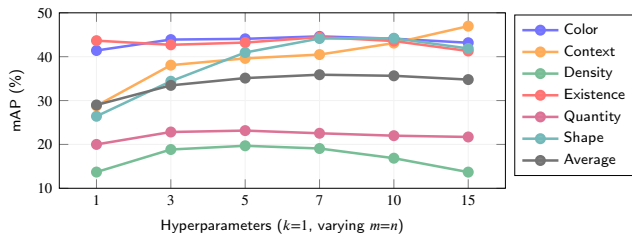
**Figure 6: Impact of modality control  $\lambda$  in WeiCom with RemoteCLIP on PatternCom.** Curves report attribute-wise mAP (%) as  $\lambda$  varies.

554 *Impact of vocabulary on FreeDM.* Figure 7 analyzes 558  
 555 how the textual memory affects FreeDM on PatternCom 589  
 556 with CLIP LAION-2B. Open Images v7 provides strong 590  
 557 performance due to its large and diverse vocabulary of 591  
 558 roughly 21k concepts. However, the EO-specific RSText vocabularies 592  
 559 remain competitive despite being much smaller, 593  
 560 with up to 200× fewer entries for RSText-150 and 20× 594  
 561 fewer entries for RSText-1k. This suggests that compact, 595  
 562 domain-grounded vocabularies can provide effective textual 596  
 563 surrogates for EO imagery. The best performance is obtained 597  
 564 by HybridText-23k, which combines Open Images with 598  
 565 RSText-2k, indicating that broad semantic coverage and EO- 599  
 566 specific grounding are complementary.



567 **Figure 7: Impact of vocabulary on FreeDM with CLIP**  
 568 LAION-2B on PatternCom. We compare Open Images v7,  
 569 EO-specific RSText vocabularies, and their hybrid combination  
 570 HybridText-23k. Bars report mAP (%) per attribute and the  
 571 average.

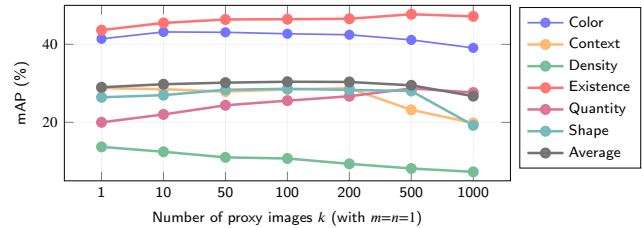
569 *Impact of textual expansion on FreeDM.* Figure 8  
 570 analyzes frequency-based textual expansion in FreeDM by  
 571 disabling visual query expansion ( $k=1$ ) and varying the  
 572 number of aggregated textual concepts ( $m=n$ ). On Pattern- 605  
 573 Com, increasing  $m=n$  improves average mAP from 28.99% 606  
 574 at  $m=n=1$  to 35.89% at  $m=n=7$ , with the largest gains 607  
 575 on shape and context. The same moderate-expansion trend  
 576 appears on xView2-CIR, where performance improves from  
 577 9.59% to 16.84% at  $m=n=7$ . In both settings, larger ex-  
 578 pansion plateau or degrade, suggesting that overly long  
 579 textual surrogates introduce noisy or off-target concepts. The  
 580 corresponding xView2-CIR analysis is reported in A.5.



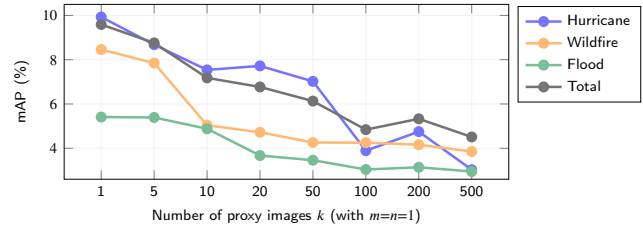
581 **Figure 8: Impact of textual expansion on FreeDM with CLIP**  
 582 LAION-2B on PatternCom. Visual expansion is disabled ( $k=1$ ),  
 583 and the number of aggregated textual concepts is varied as  
 584  $m=n$ .

584 *Impact of visual expansion on FreeDM.* Figure 9  
 585 and 10 analyze visual query expansion in FreeDM by  
 586 varying the number of proxy images  $k$  while disabling  
 587 textual expansion ( $m=n=1$ ). On PatternCom, visual

588 expansion is *beneficial* up to a moderate range: average  
 589 mAP increases from 28.99% at  $k=1$  to 30.41% at  $k=100$ ,  
 590 before dropping to 26.71% at  $k=1000$ . This is consistent  
 591 with the *same class + target attribute* criterion, where  
 592 nearest neighbors often remain class-consistent and provide  
 593 useful attribute evidence, until overly large proxy sets  
 594 introduce noisy distractors. In contrast, visual expansion is  
 595 *harmful* on xView2-CIR: Total mAP decreases from 9.59%  
 596 at  $k=1$  to 4.51% at  $k=500$ . Since xView2-CIR requires  
 597 the *same scene/location + target state*, visually similar  
 598 proxies usually correspond to different locations and inject  
 599 irrelevant content. This explains why proxy-based expansion  
 600 benefits attribute-oriented retrieval but can degrade identity-  
 601 preserving change retrieval.



602 **Figure 9: Impact of visual expansion on FreeDM with CLIP**  
 603 LAION-2B on PatternCom. The number of proxy images  $k$   
 604 is varied while textual expansion is disabled ( $m=n=1$ ).



605 **Figure 10: Impact of visual expansion on FreeDM with**  
 606 OpenAI CLIP on xView2-CIR. The number of proxy images  
 607  $k$  is varied while textual expansion is disabled ( $m=n=1$ ).

608 *Ablation of BASIC components.* We further analyze  
 609 the contribution of BASIC components, with full results  
 610 reported in A.5. On PatternCom, the largest drops occur  
 611 when removing *centering* or *semantic projection*, indicating  
 612 that modality calibration and semantic subspace alignment  
 613 are the main drivers of BASIC’s performance. Harris  
 614 weighting, min-based normalization, text contextualization,  
 615 and query expansion provide smaller, secondary gains. A  
 616 similar trend appears on xView2-CIR, where centering and  
 617 semantic projection remain the most important components,  
 618 while visual query expansion becomes harmful under the  
 619 *same scene/location + target state* criterion. This  
 620 supports the broader finding that proxy-based expansion can  
 621 help attribute-oriented retrieval but may degrade identity-  
 622 preserving change retrieval.

623 *Sensitivity to text query phrasing.* We also evaluate  
 624 whether composed retrieval is sensitive to the phrasing of  
 625 the textual modifier, with full results reported in A.5. On  
 626 PatternCom, the original plain attribute values (e.g., blue,  
 627 water, dense) are generally the most reliable, especially for

stronger multimodal methods. More verbose templates such as having blue color or with the main object modified to have blue color often reduce performance, suggesting that additional function words can dilute the compact attribute signal in this benchmark.

On xView2-CIR, prompt sensitivity is more dataset- and method-dependent. Rephrasing post-`{disaster}` into impact-oriented or disaster-explicit descriptions (e.g., burned area, flooded area, fire-affected region) can benefit methods that rely more strongly on semantic text structure, particularly BASIC. Overall, these results indicate that prompt phrasing is an important practical factor: compact modifiers are preferable for controlled attribute edits, whereas change-centric retrieval may benefit from more descriptive event semantics.

#### 4.5. Discussion

*Generalization across backbones and EO adaptation.* Our benchmark suggests that both remote-sensing adaptation and general VLM quality matter for RSCIR. RS-adapted CLIP variants, such as RemoteCLIP, CLIP LAION-RS, and SkyCLIP-50, often improve over generic CLIP-style backbones under the same protocol. However, the gains from domain adaptation can be comparable to, or smaller than, the gains obtained by using stronger general-purpose backbones and training objectives, such as SigLIP. This indicates that advances in VLM architecture, pre-training data, and loss design transfer directly to composed retrieval in EO imagery. Consequently, future RSCIR systems should consider both EO-specific adaptation and improvements in general VLM representation quality.

*Scalability and deployment trade-offs.* The results show that training-free composition methods provide strong baselines for EO retrieval. FreeDom is the strongest method on PatternCom, while WeiCom and BASIC remain competitive lightweight alternatives. This is important for operational EO systems, where latency, scalability, and engineering complexity matter alongside accuracy. CIReVL can be effective, but requires image captioning and LLM-based rewriting, introducing additional computational cost and possible failure modes. In contrast, score-fusion or feature-calibration methods such as WeiCom and BASIC are simpler to deploy over large archives. Supervised multimodal heads such as MagicLens also transfer to EO imagery to some extent, but their performance may remain limited without training on remote-sensing compositional data.

*Attribute-based and change-centric RSCIR require different mechanisms.* The contrast between PatternCom and xView2-CIR shows that change-centric RSCIR should be treated as a distinct retrieval setting rather than a direct extension of attribute-based retrieval. PatternCom rewards methods that preserve class identity while injecting attribute-specific textual cues, which explains the strong performance of vocabulary-based methods such as FreeDom. In xView2-CIR, however, relevance requires the same scene/location and a target post-event state. This makes

proxy-based visual expansion less reliable, because visually similar images often correspond to different locations and can dilute the scene-identity signal. These findings suggest that methods designed for attribute editing may not directly transfer to identity-preserving change retrieval.

*Role in operational EO workflows.* RSCIR is not a replacement for pixel-level change detection, damage assessment, or rapid mapping. Instead, it provides a complementary interface for controllable archive exploration and rapid evidence gathering. In practical workflows, composed retrieval can support location-aware search, retrieval of visually similar historical cases, analyst-in-the-loop exploration, quality control through hard-negative discovery, and weak supervision for downstream change models. This is especially useful when dense annotations, perfect co-registration, or complete temporal coverage are unavailable.

*Interpretability and analyst interaction.* Interpretability is also important for EO deployment. Methods based on explicit textual memories, such as FreeDom, provide human-readable surrogate concepts for the visual query, making it easier to inspect which semantic cues drive retrieval. Score-fusion methods such as WeiCom expose the relative contribution of image and text through the modality weight, while BASIC highlights the role of feature calibration and semantic projection. Such properties can help analysts diagnose whether failures arise from excessive visual dominance, weak text grounding, or scene-identity drift.

*Sensitivity to text phrasing.* Prompt phrasing is a practical factor in composed retrieval. On PatternCom, compact attribute modifiers such as `blue`, `water`, or `dense` are generally most effective, whereas verbose templates can dilute the attribute signal. In contrast, xView2-CIR can benefit from more descriptive, impact-oriented formulations, particularly for methods that rely more strongly on semantic text structure. This suggests that prompt design should be treated as part of the retrieval system: short modifiers are suitable for controlled attribute edits, while change-centric settings may require more explicit event or impact descriptions.

## 5. Limitations

This study has several limitations. First, although xView2-CIR extends RSCIR toward operationally relevant change-centric retrieval, it remains relatively small and imbalanced, especially for rare disaster categories; results on these categories should therefore be interpreted as indicative rather than definitive. Second, our evaluation focuses on image-level retrieval and does not address finer-grained localization or dense temporal search. Third, some composition methods are more naturally tied to specific backbones or representation spaces, which limits strict cross-backbone comparability. Finally, prompt design and auxiliary vocabularies influence performance, particularly in change-centric scenarios, and require more systematic

study. We therefore view this work as a benchmark-and-analysis foundation for RSCIR rather than a final solution to composed retrieval in Earth observation.

## 6. Conclusion

We studied remote sensing composed image retrieval (RSCIR), where a query combines a reference image with a textual modifier to express targeted retrieval intent. We established a unified benchmark on PatternCom with domain-grounded adaptations of representative composition methods and a standardized protocol spanning six vision-language backbones. To connect benchmarking with operational EO needs, we introduced xView2-CIR, a change-centric benchmark for disaster and damage monitoring, where retrieval requires both scene identity and a target post-event state.

Our experiments show that stronger vision-language backbones transfer directly to RSCIR, and that training-free composition strategies can be highly effective in EO settings. FreeDom achieves the strongest performance on PatternCom, while lightweight methods such as WeiCom and BASIC remain competitive and attractive for practical deployment. At the same time, xView2-CIR reveals that change-centric retrieval differs substantially from attribute-based editing: preserving scene identity changes the method ranking and exposes weaknesses in proxy-based expansion strategies.

Overall, our results position RSCIR as a practical complement to standard RSIR and change-analysis pipelines, especially when analysts need controllable, semantically guided access to large EO archives. Future work should expand change-centric benchmarks, study prompt and vocabulary design more systematically, and connect scene-level composed retrieval with finer-grained localization and temporal reasoning.

## 7. Data and Code Availability Statement

The code and data supporting the findings of this study are publicly available at <https://github.com/billpsomas/rscir>. PatternCom is based on publicly available source data. The proposed xView2-CIR benchmark and the scripts used for data preparation, evaluation, and reproduction of the reported results are also provided in the repository.

## 8. Acknowledgment

This work was supported by the Czech Technical University in Prague grant No. SGS23/173/OHK3/3T/13, the EU Horizon Europe programme MSCA PF RAVIOLI (No. 101205297), and the Junior Star GACR GM 21-28830M. We acknowledge VSB-Technical University of Ostrava, IT4Innovations National Supercomputing Center, Czech Republic, for awarding this project (OPEN-33-67) access to the LUMI supercomputer, owned by the EuroHPC Joint Undertaking, hosted by CSC (Finland) and the LUMI consortium, through the Ministry of

Education, Youth and Sports of the Czech Republic via the e-INFRA CZ project (ID: 90254). The access to the computational infrastructure of the OP VVV funded project CZ.02.1.01/0.0/0.0/16\_019/0000765 “Research Center for Informatics” is also gratefully acknowledged.

## References

- [1] Agouris, P., Carswell, J., Stefanidis, A., 1999. An environment for content-based image retrieval from large spatial databases. *ISPRS J. Photogramm. Remote Sens.*
- [2] Baldradi, A., Agnolucci, L., Bertini, M., Del Bimbo, A., 2023. Zero-shot composed image retrieval with textual inversion, in: *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, pp. 15338–15347.
- [3] Baldradi, A., Bertini, M., Uricchio, T., Del Bimbo, A., 2022. Effective conditioned and composed image retrieval combining clip-based features, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*
- [4] Berg, T.L., Berg, A.C., Shih, J., 2010. Automatic attribute discovery and characterization from noisy web data, in: *Proc. Eur. Conf. Comput. Vis.*, Springer.
- [5] Bhagavathy, S., Manjunath, B.S., 2006. Modeling and detection of geospatial objects using texture motifs. *IEEE Trans. Geosci. Remote Sens.* 44, 3706–3715.
- [6] Boualleg, Y., Farah, M., 2018. Enhanced interactive remote sensing image retrieval with scene classification convolutional neural networks model, in: *Proc. IEEE Int. Geosci. Remote Sens. Symp.*, pp. 4748–4751.
- [7] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al., 2020. Language models are few-shot learners. *Adv. Neural Inf. Process. Syst.* 33, 1877–1901.
- [8] Cao, R., Zhang, Q., Zhu, J., Li, Q., Li, Q., Liu, B., Qiu, G., 2020. Enhancing remote sensing image retrieval using a triplet deep metric learning network. *Int. J. Remote Sens.* 41, 740–751.
- [9] Chaudhuri, B., Demir, B., Chaudhuri, S., Bruzzone, L., 2017. Multilabel remote sensing image retrieval using a semisupervised graph-theoretic method. *IEEE Trans. Geosci. Remote Sens.* 56, 1144–1158.
- [10] Chaudhuri, U., Banerjee, B., Bhattacharya, A., 2019. Siamese graph convolutional network for content based remote sensing image retrieval. *Comput. Vis. Image Underst.*
- [11] Chaudhuri, U., Banerjee, B., Bhattacharya, A., Datcu, M., 2021. Attention-driven graph convolution network for remote sensing image retrieval. *IEEE Geosci. Remote Sens. Lett.* 19, 1–5.
- [12] Chen, Y., Bazzani, L., 2020. Learning joint visual semantic matching embeddings for language-guided retrieval, in: *Proc. Eur. Conf. Comput. Vis.*
- [13] Chen, Y., Gong, S., Bazzani, L., 2020. Image search with text feedback by visiolinguistic attention learning, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*
- [14] Cheng, Q., Gan, D., Fu, P., Huang, H., Zhou, Y., 2021a. A novel ensemble architecture of residual attention-based deep metric learning for remote sensing image retrieval. *Remote Sens.* 13, 3445.
- [15] Cheng, Q., Huang, H., Ye, L., Fu, P., Gan, D., Zhou, Y., 2021b. A semantic-preserving deep hashing model for multi-label remote sensing image retrieval. *Remote Sens.* 13, 4965.
- [16] Cherti, M., Beaumont, R., Wightman, R., Zhai, X., Beyer, L., Kolesnikov, A., Dosovitskiy, A., Houthby, N., Minderer, M., 2022. Openclip: An open source implementation of clip.
- [17] Dai, O.E., Demir, B., Sankur, B., Bruzzone, L., 2017. A novel system for content based retrieval of multi-label remote sensing images, in: *Proc. IEEE Int. Geosci. Remote Sens. Symp.*, pp. 1744–1747.
- [18] Delmas, G., de Rezende, R.S., Csurka, G., Larlus, D., 2022. Artemis: Attention-based retrieval with text-explicit matching and implicit similarity.
- [19] Efthymiadis, N., Psomas, B., Laskar, Z., Karantzas, K., Avrithis, Y., Chum, O., Tolia, G., 2025. Composed image retrieval for

852 training-free domain conversion, in: Proc. IEEE/CVF Winter Conf. 920  
 Appl. Comput. Vis. 921

853 [20] Fan, L., Zhao, H., Zhao, H., 2020. Global optimization: Combining 922  
 854 local loss with result ranking loss in remote sensing image retrieval. 923  
 855 IEEE Trans. Geosci. Remote Sens. 59, 7011–7026. 924

856 [21] Frome, A., Corrado, G.S., Shlens, J., Bengio, S., Dean, J., Ranzato, 925  
 857 M., Mikolov, T., 2013. Devise: A deep visual-semantic embedding 926  
 858 model. Adv. Neural Inf. Process. Syst. 26. 927

859 [22] Ge, Y., Jiang, S., Xu, Q., Jiang, C., Ye, F., 2018. Exploiting 928  
 860 representations from pre-trained convolutional neural networks for 929  
 861 high-resolution remote sensing image retrieval. Multimed. Tools 930  
 862 Appl. 77, 17489–17515. 931

863 [23] Gordo, A., Almazan, J., Revaud, J., Larlus, D., 2016. Deep image 932  
 864 retrieval: Learning global representations for image search, in: Proc. 933  
 865 Eur. Conf. Comput. Vis. 934

866 [24] Gorelick, N., Hancher, M., Dixon, M., Ilyushchenko, S., Thau, D., 935  
 867 Moore, R., 2017. Google earth engine: Planetary-scale geospatial 936  
 868 analysis for everyone. Remote Sens. Environ. 202, 18–27. 937

869 [25] Gu, G., Chun, S., Kim, W., Jun, H., Kang, Y., Yun, S., a. Compodiff: 938  
 870 Versatile composed image retrieval with latent diffusion. Transactions 939  
 871 on Machine Learning Research . 940

872 [26] Gu, X., Lin, T.Y., Kuo, W., Cui, Y., b. Open-vocabulary object 941  
 873 detection via vision and language knowledge distillation. 942

874 [27] Haklay, M., Weber, P., 2008. Openstreetmap: User-generated street 943  
 875 maps. IEEE Pervasive Comput. 7, 12–18. 944

876 [28] Han, X., Wu, Z., Huang, P.X., Zhang, X., Zhu, M., Li, Y., Zhao, 945  
 877 Y., Davis, L.S., 2017. Automatic spatially-aware fashion concept 946  
 878 discovery, in: Proc. IEEE/CVF Int. Conf. Comput. Vis. 947

879 [29] Hosseinzadeh, M., Wang, Y., 2020. Composed query image retrieval 948  
 880 using locally bounded features, in: Proc. IEEE/CVF Conf. Comput. 949  
 881 Vis. Pattern Recognit. 950

882 [30] Hou, D., Miao, Z., Xing, H., Wu, H., 2020. Exploiting low 951  
 883 dimensional features from the mobilenets for remote sensing image 952  
 884 retrieval. Earth Sci. Inform. 13, 1437–1443. 953

885 [31] Hu, F., Tong, X., Xia, G.S., Zhang, L., 2016. Delving into deep 954  
 886 representations for remote sensing image retrieval, in: Proc. IEEE 955  
 887 Int. Conf. Signal Process., pp. 198–203. 956

888 [32] Hu, S., Feng, M., Nguyen, R.M., Lee, G.H., 2018. Cvm-net: 957  
 889 Cross-view matching network for image-based ground-to-aerial geo- 958  
 890 localization, in: Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 959  
 891 pp. 7258–7267. 960

892 [33] Imbriaco, R., Sebastian, C., Bondarev, E., de With, P.H., 2019. 961  
 893 Aggregated deep local features for remote sensing image retrieval. 962  
 894 Remote Sens. 11, 493. 963

895 [34] Imbriaco, R., Sebastian, C., Bondarev, E., de With, P.H., 2021. 964  
 896 Toward multilabel image retrieval for remote sensing. IEEE Trans. 965  
 897 Geosci. Remote Sens. 60, 1–14. 966

898 [35] Isola, P., Lim, J.J., Adelson, E.H., 2015. Discovering states and 967  
 899 transformations in image collections, in: Proc. IEEE/CVF Conf. 968  
 900 Comput. Vis. Pattern Recognit., pp. 1383–1391. 969

901 [36] Jia, C., Yang, Y., Xia, Y., Chen, Y.T., Parekh, Z., Pham, H., Le, 970  
 902 Q., Sung, Y.H., Li, Z., Duerig, T., 2021. Scaling up visual and 971  
 903 vision-language representation learning with noisy text supervision, 972  
 904 in: Proc. Int. Conf. Mach. Learn. 973

905 [37] Kang, J., Fernandez-Beltran, R., Hong, D., Chanussot, J., Plaza, A., 974  
 906 2020. Graph relation network: Modeling relations between scenes 975  
 907 for multilabel remote-sensing image classification and retrieval. 976  
 908 IEEE Trans. Geosci. Remote Sens. 59, 4355–4369. 977

909 [38] Karthik, S., Roth, K., Mancini, M., Akata, Z., 2024. Vision-by- 978  
 910 language for training-free compositional image retrieval, in: Int. 979  
 911 Conf. Learn. Represent., pp. 16926–16941. 980

912 [39] Khurshid, N., Hanif, T., Tharani, M., Taj, M., 2019. Cross-view im- 981  
 913 age retrieval-ground to aerial image retrieval through deep learning, 982  
 914 in: Proc. Int. Conf. Neural Inf. Process., Springer. pp. 210–221. 983

915 [40] Kuznetsova, A., Rom, H., Alldrin, N., Uijlings, J., Krasin, I., Pont- 984  
 916 Tuset, J., Kamali, S., Popov, S., Mallocci, M., Kolesnikov, A., et al., 985  
 917 2020. The open images dataset v4: Unified image classification, 986  
 918 object detection, and visual relationship detection at scale. Int. J. 987  
 Comput. Vis. .

[41] Lam, D., Kuzma, R., McGee, K., Dooley, S., Laielli, M., Klaric, 921  
 M., Bulatov, Y., McCord, B., 2018. xview: Objects in context in 922  
 overhead imagery. arXiv preprint arXiv:1802.07856 . 923

[42] Lee, S., Kim, D., Han, B., 2021. Cosmo: Content-style modulation 924  
 for image retrieval with text feedback, in: Proc. IEEE/CVF Conf. 925  
 Comput. Vis. Pattern Recognit. 926

[43] Li, J., Li, D., Savarese, S., Hoi, S., 2023. Blip-2: Bootstrapping 927  
 language-image pre-training with frozen image encoders and large 928  
 language models, in: Proc. Int. Conf. Mach. Learn., pp. 19730– 929  
 19742. 930

[44] Li, J., Li, D., Xiong, C., Hoi, S., 2022. Blip: Bootstrapping language- 931  
 image pre-training for unified vision-language understanding and 932  
 generation, in: Proc. Int. Conf. Mach. Learn. 933

[45] Li, J., Narayanan, R.M., 2004. Integrated spectral and spatial 934  
 information mining in remote sensing imagery. IEEE Trans. Geosci. 935  
 Remote Sens. 42, 673–685. 936

[46] Li, Y., Zhang, Y., Huang, X., Ma, J., 2018. Learning source-invariant 937  
 deep hashing convolutional neural networks for cross-source remote 938  
 sensing image retrieval. IEEE Trans. Geosci. Remote Sens. 56, 939  
 6521–6536. 940

[47] Li, Y., Zhang, Y., Tao, C., Zhu, H., 2016. Content-based high- 941  
 resolution remote sensing image retrieval via unsupervised feature 942  
 learning and collaborative affinity metric fusion. Remote Sens. 8, 943  
 709. 944

[48] Lin, T.Y., Cui, Y., Belongie, S., Hays, J., 2015. Learning deep 945  
 representations for ground-to-aerial geolocalization, in: Proc. IEEE 946  
 Conf. Comput. Vis. Pattern Recognit., pp. 5007–5015. 947

[49] Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, 948  
 D., Dollár, P., Zitnick, C.L., 2014. Microsoft coco: Common objects 949  
 in context, in: Proc. Eur. Conf. Comput. Vis. 950

[50] Liu, C., Zhao, R., Chen, H., Zou, Z., Shi, Z., 2022. Remote sensing 951  
 image change captioning with dual-branch transformers: A new 952  
 method and a large scale dataset. IEEE Trans. Geosci. Remote Sens. 953  
 60, 1–20. 954

[51] Liu, F., Chen, D., Guan, Z., Zhou, X., Zhu, J., Ye, Q., Fu, L., Zhou, J., 955  
 2024. Remoteclip: A vision language foundation model for remote 956  
 sensing. IEEE Trans. Geosci. Remote Sens. 62, 1–16. 957

[52] Liu, Y., Ding, L., Chen, C., Liu, Y., 2020a. Similarity-based unsu- 958  
 pervised deep transfer learning for remote sensing image retrieval. 959  
 IEEE Trans. Geosci. Remote Sens. 58, 7872–7889. 960

[53] Liu, Y., Han, Z., Chen, C., Ding, L., Liu, Y., 2020b. Eagle-eyed 961  
 multitask cnns for aerial image retrieval and scene classification. 962  
 IEEE Trans. Geosci. Remote Sens. 58, 6699–6721. 963

[54] Liu, Y., Liu, Y., Chen, C., Ding, L., 2020c. Remote-sensing image 964  
 retrieval with tree-triplet-classification networks. Neurocomputing 965  
 405, 48–61. 966

[55] Lv, Y., Xiong, W., Zhang, X., Cui, Y., 2021. Fusion-based correla- 967  
 tion learning model for cross-modal remote sensing image retrieval. 968  
 IEEE Geosci. Remote Sens. Lett. 19, 1–5. 969

[56] Ma, C., Dai, Q., Liu, J., Liu, S., Yang, J., 2014. An improved svm 970  
 model for relevance feedback in remote sensing image retrieval. Int. 971  
 J. Digit. Earth 7, 725–745. 972

[57] Ma, J., Shi, D., Tang, X., Zhang, X., Han, X., Jiao, L., 2021. Cross- 973  
 source image retrieval based on ensemble learning and knowledge 974  
 distillation for remote sensing images, in: Proc. IEEE Int. Geosci. 975  
 Remote Sens. Symp., IEEE. pp. 2803–2806. 976

[58] Mamatha, Y., Ananth, A., 2010. Content based image retrieval 977  
 of satellite imageries using soft query based color composite tech- 978  
 niques. Int. J. Comput. Appl. 7, 0975–8887. 979

[59] Napoletano, P., 2018. Visual descriptors for content-based retrieval 980  
 of remote-sensing images. Int. J. Remote Sens. 39, 1343–1376. 981

[60] Neculai, A., Chen, Y., Akata, Z., 2022. Probabilistic compositional 982  
 embeddings for multimodal image retrieval, in: Proc. IEEE/CVF 983  
 Conf. Comput. Vis. Pattern Recognit. 984

[61] Noh, H., Araujo, A., Sim, J., Weyand, T., Han, B., 2017. Large- 985  
 scale image retrieval with attentive deep local features, in: Proc. 986  
 IEEE/CVF Int. Conf. Comput. Vis. 987

- 988 [62] OpenAI, 2024. Gpt-4o system card. 1055
- 989 [63] Piedra-Fernandez, J.A., Ortega, G., Wang, J.Z., Canton-Garbin, M., 1056  
990 2013. Fuzzy content-based image retrieval for oceanic remote 1057  
991 sensing. *IEEE Trans. Geosci. Remote Sens.* 52, 5422–5431. 1058
- 992 [64] Psomas, B., Kakogeorgiou, I., Efthymiadis, N., Toliás, G., Chum, 1059  
993 O., Avrithis, Y., Karantzalos, K., 2024. Composed image retrieval 1060  
994 for remote sensing, in: *Proc. IEEE Int. Geosci. Remote Sens. Symp.* 1061
- 995 [65] Psomas, B., Retsinas, G., Efthymiadis, N., Filntisis, P., Avrithis, Y., 1062  
996 Maragos, P., Chum, O., Toliás, G., 2025. Instance-level composed 1063  
997 image retrieval, in: *Adv. Neural Inf. Process. Syst.* 1064
- 998 [66] Radenović, F., Toliás, G., Chum, O., 2019. Fine-tuning cnn image 1065  
999 retrieval with no human annotation. *IEEE Trans. Pattern Anal.* 1066  
1000 *Mach. Intell.* . 1067
- 1001 [67] Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, 1068  
1002 S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al., 2021. Learning 1069  
1003 transferable visual models from natural language supervision, in: 1070  
1004 *Proc. Int. Conf. Mach. Learn.* 1071
- 1005 [68] Ren, S., He, K., Girshick, R., Sun, J., 2015. Faster r-cnn: Towards 1072  
1006 real-time object detection with region proposal networks. *Adv.* 1073  
1007 *Neural Inf. Process. Syst.* . 1074
- 1008 [69] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B., 2022. 1075  
1009 High-resolution image synthesis with latent diffusion models, in: 1076  
1010 *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.* 1077
- 1011 [70] Sadeghi-Tehran, P., Angelov, P., Viret, N., Hawkesford, M.J., 2019. 1078  
1012 Scalable database indexing and fast image retrieval based on deep 1079  
1013 learning and hierarchically nested structure applied to remote sens- 1080  
1014 ing and plant biology. *J. Imaging* 5, 33. 1081
- 1015 [71] Saito, K., Sohn, K., Zhang, X., Li, C.L., Lee, C.Y., Saenko, K., 1082  
1016 Pfister, T., 2023. Pic2word: Mapping pictures to words for zero-shot 1083  
1017 composed image retrieval, in: *Proc. IEEE/CVF Conf. Comput. Vis.* 1084  
1018 *Pattern Recognit.* 1085
- 1019 [72] Sarafianos, N., Xu, X., Kakadiaris, I.A., 2019. Adversarial represen- 1086  
1020 tation learning for text-to-image matching, in: *Proc. IEEE/CVF Int.* 1087  
1021 *Conf. Comput. Vis.*, pp. 5814–5824. 1088
- 1022 [73] Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, 1089  
1023 R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., 1090  
1024 et al., 2022. Laion-5b: An open large-scale dataset for training next 1091  
1025 generation image-text models. *Adv. Neural Inf. Process. Syst.* . 1092
- 1026 [74] Shao, Z., Zhou, W., Cheng, Q., 2014. Remote sensing image retrieval 1093  
1027 with combined features of salient region. *Int. Arch. Photogramm.* 1094  
1028 *Remote Sens. Spatial Inf. Sci.* 40, 83–88. 1095
- 1029 [75] Shao, Z., Zhou, W., Deng, X., Zhang, M., Cheng, Q., 2020. Mul- 1096  
1030 tilabel remote sensing image retrieval based on fully convolutional 1097  
1031 network. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.* 13, 318– 1098  
1032 328. 1099
- 1033 [76] Shi, Y., Li, H., 2022. Beyond cross-view image retrieval: Highly ac- 1100  
1034 curate vehicle localization using satellite image, in: *Proc. IEEE/CVF* 1101  
1035 *Conf. Comput. Vis. Pattern Recognit.*, pp. 17010–17020. 1102
- 1036 [77] Stojnić, V., Kalantidis, Y., Matas, J., Toliás, G., 2025. Lploss: Label 1103  
1037 propagation over patches and pixels for open-vocabulary seman- 1104  
1038 tic segmentation, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern* 1105  
1039 *Recognit.*, pp. 9794–9803. 1106
- 1040 [78] Sumbul, G., De Wall, A., Kreuziger, T., Marcelino, F., Costa, 1107  
1041 H., Benevides, P., Caetano, M., Demir, B., Markl, V., 2021a. 1108  
1042 Bigearthnet-mm: A large-scale, multimodal, multilabel benchmark 1109  
1043 archive for remote sensing image classification and retrieval. *IEEE* 1110  
1044 *Geosci. Remote Sens. Mag.* 9, 174–180. 1111
- 1045 [79] Sumbul, G., Demir, B., 2021. A novel graph-theoretic deep rep- 1112  
1046 resentation learning method for multi-label remote sensing image 1113  
1047 retrieval, in: *Proc. IEEE Int. Geosci. Remote Sens. Symp.*, IEEE. 1114  
1048 pp. 266–269. 1115
- 1049 [80] Sumbul, G., Demir, B., 2022. Plasticity-stability preserving multi- 1116  
1050 task learning for remote sensing image retrieval. *IEEE Trans.* 1117  
1051 *Geosci. Remote Sens.* 60, 1–16. 1118
- 1052 [81] Sumbul, G., Ravanbakhsh, M., Demir, B., 2021b. Informative and 1119  
1053 representative triplet selection for multilabel remote sensing image 1120  
1054 retrieval. *IEEE Trans. Geosci. Remote Sens.* 60, 1–11. 1121
- [82] Sun, Y., Feng, S., Ye, Y., Li, X., Kang, J., Huang, Z., Luo, C., 2021. Multisensor fusion and explicit semantic preserving-based deep hashing for cross-modal remote sensing image retrieval. *IEEE Trans. Geosci. Remote Sens.* 60, 1–14.
- [83] Tang, X., Zhang, X., Liu, F., Jiao, L., 2018. Unsupervised deep feature learning for remote sensing image retrieval. *Remote Sens.* 10, 1243.
- [84] Tewel, Y., Shalev, Y., Schwartz, I., Wolf, L., 2022. Zerocap: Zero-shot image-to-text generation for visual-semantic arithmetic, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pp. 17918–17928.
- [85] Tian, Y., Deng, X., Zhu, Y., Newsam, S., 2020. Cross-time and orientation-invariant overhead image geolocation using deep local features, in: *Proc. IEEE/CVF Winter Conf. Appl. Comput. Vis.*, pp. 2512–2520.
- [86] Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al., 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288* .
- [87] Vo, N., Jiang, L., Sun, C., Murphy, K., Li, L.J., Fei-Fei, L., Hays, J., 2019. Composing text and image for image retrieval-an empirical odyssey, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*
- [88] Wang, M., Song, T., 2012. Remote sensing image retrieval by scene semantic matching. *IEEE Trans. Geosci. Remote Sens.* 51, 2874–2886.
- [89] Wang, M., Wan, Q., Gu, L., Song, T., 2013. Remote-sensing image retrieval by combining image visual and semantic features. *Int. J. Remote Sens.* 34, 4200–4223.
- [90] Wang, S., Hou, D., Xing, H., 2022. A novel multi-attention fusion network with dilated convolution and label smoothing for remote sensing image retrieval. *Int. J. Remote Sens.* 43, 1306–1322.
- [91] Wang, Y., Ji, S., Lu, M., Zhang, Y., 2020. Attention boosted bilinear pooling for remote sensing image retrieval. *Int. J. Remote Sens.* 41, 2704–2724.
- [92] Wang, Y., Ji, S., Zhang, Y., 2021. A learnable joint spatial and spectral transformation for high resolution remote sensing image retrieval. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.* 14, 8100–8112.
- [93] Wang, Y., Zhang, L., Tong, X., Zhang, L., Zhang, Z., Liu, H., Xing, X., Mathiopoulos, P.T., 2016. A three-layered graph-based learning approach for remote sensing image retrieval. *IEEE Trans. Geosci. Remote Sens.* 54, 6020–6034.
- [94] Wang, Z., Prabha, R., Huang, T., Wu, J., Rajagopal, R., 2024. Skyscript: A large and semantically diverse vision-language dataset for remote sensing, in: *Proc. AAAI Conf. Artif. Intell.*, pp. 5805–5813.
- [95] Wu, H., Gao, Y., Guo, X., Al-Halah, Z., Rennie, S., Grauman, K., Feris, R., 2021. Fashion iq: A new dataset towards retrieving images by natural language feedback, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*
- [96] Xiong, W., Lv, Y., Cui, Y., Zhang, X., Gu, X., 2019. A discriminative feature learning approach for remote sensing image retrieval. *Remote Sens.* 11, 281.
- [97] Xiong, W., Lv, Y., Zhang, X., Cui, Y., 2020a. Learning to translate for cross-source remote sensing image retrieval. *IEEE Trans. Geosci. Remote Sens.* 58, 4860–4874.
- [98] Xiong, W., Xiong, Z., Cui, Y., Lv, Y., 2020b. A discriminative distillation network for cross-source remote sensing image retrieval. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.* 13, 1234–1247.
- [99] Xu, F., Yang, W., Jiang, T., Lin, S., Luo, H., Xia, G.S., 2020. Mental retrieval of remote sensing images via adversarial sketch-image feature learning. *IEEE Trans. Geosci. Remote Sens.* 58, 7801–7814.
- [100] Ye, F., Xiao, H., Zhao, X., Dong, M., Luo, W., Min, W., 2018. Remote sensing image retrieval using convolutional neural network features and weighted distance. *IEEE Geosci. Remote Sens. Lett.* 15, 1535–1539.

- 1122 [101] Yin, M., Yao, Z., Cao, Y., Li, X., Zhang, Z., Lin, S., Hu, H.,  
1123 2020. Disentangled non-local neural networks, in: Proc. Eur. Conf.  
1124 Comput. Vis.
- 1125 [102] Yuan, Z., Zhang, W., Fu, K., Li, X., Deng, C., Wang, H., Sun, X.,  
1126 2021. Exploring a fine-grained multiscale method for cross-modal  
1127 remote sensing image retrieval. *IEEE Trans. Geosci. Remote Sens.*  
1128 60, 1–19.
- 1129 [103] Yuan, Z., Zhang, W., Tian, C., Rong, X., Zhang, Z., Wang, H., Fu,  
1130 K., Sun, X., 2022. Remote sensing cross-modal text-image retrieval  
1131 based on global and local information. *IEEE Trans. Geosci. Remote*  
1132 *Sens.* 60, 1–16.
- 1133 [104] Zeng, Z., Wang, Z., Yang, F., Satoh, S., 2022. Geo-localization  
1134 via ground-to-satellite cross-view image retrieval. *IEEE Trans.*  
1135 *Multimedia* .
- 1136 [105] Zhai, X., Mustafa, B., Kolesnikov, A., Beyer, L., 2023. Sigmoid  
1137 loss for language image pre-training, in: Proc. IEEE/CVF Int. Conf.  
1138 Comput. Vis., pp. 11975–11986.
- 1139 [106] Zhang, K., Luan, Y., Hu, H., Lee, K., Qiao, S., Chen, W., Su, Y.,  
1140 Chang, M.W., 2024. Magiclens: Self-supervised image retrieval  
1141 with open-ended instructions, in: Proc. Int. Conf. Mach. Learn., pp.  
1142 59403–59420.
- 1143 [107] Zhang, M., Cheng, Q., Luo, F., Ye, L., 2021. A triplet nonlocal neural  
1144 network with dual-anchor triplet loss for high-resolution remote  
1145 sensing image retrieval. *IEEE J. Sel. Top. Appl. Earth Obs. Remote*  
1146 *Sens.* 14, 2711–2723.
- 1147 [108] Zhang, Q., Lei, Z., Zhang, Z., Li, S.Z., 2020. Context-aware attention  
1148 network for image-text retrieval, in: Proc. IEEE/CVF Conf. Comput.  
1149 Vis. Pattern Recognit., pp. 3536–3545.
- 1150 [109] Zhao, H., Yuan, L., Zhao, H., Wang, Z., 2021. Global-aware ranking  
1151 deep metric learning for remote sensing image retrieval. *IEEE*  
1152 *Geosci. Remote Sens. Lett.* 19, 1–5.
- 1153 [110] Zhou, W., Guan, H., Li, Z., Shao, Z., Delavar, M.R., 2023. Remote  
1154 sensing image retrieval in the past decade: Achievements, chal-  
1155 lenges, and future directions. *IEEE J. Sel. Top. Appl. Earth Obs.*  
1156 *Remote Sens.* .
- 1157 [111] Zhou, W., Newsam, S., Li, C., Shao, Z., 2017. Learning low di-  
1158 mensional convolutional neural networks for high-resolution remote  
1159 sensing image retrieval. *Remote Sens.* 9, 489.
- 1160 [112] Zhou, W., Newsam, S., Li, C., Shao, Z., 2018. Patternnet: A  
1161 benchmark dataset for performance evaluation of remote sensing  
1162 image retrieval. *ISPRS J. Photogramm. Remote Sens.* 145, 197–209.
- 1163 [113] Zhuo, Z., Zhou, Z., 2021. Remote sensing image retrieval with  
1164 gabor-ca-resnet and split-based deep feature transform network.  
1165 *Remote Sens.* 13, 869.

1166 **Appendix**

1167 **A. Additional Experiments**

1168 **A.1. Dataset statistics**

1169 **Table 4** reports the detailed statistics of PatternCom, including the number of composed queries and positives for each attribute type, class, and target value. These statistics highlight the variability in the number of positives across attribute values, motivating the use of attribute-balanced macro-averaged mAP in the main evaluation.

**Table 4**

**PatternCom statistics.** Breakdown by attribute type, class, and attribute value. #Queries denotes the number of composed queries where the query text specifies the corresponding target value, and #Positives denotes the number of database images that satisfy the relevance criterion (*same class and target attribute value*).

ATTRIBUTE	CLASS	VALUE	#QUERIES	#POSITIVES	
color	airplane	white	53	672	
		purple	672	53	
	nursing home	white	383	85	
		gray	85	383	
	crosswalk	white	388	412	
		yellow	412	388	
	tennis court	blue	287	339	
		brown	624	2	
		gray	576	50	
		green	415	211	
red		602	24		
context	bridge	concrete	800	800	
		water	800	800	
density	residential	sparse	800	800	
		dense	800	800	
existence	parking	with cars	653	947	
		without cars	947	653	
	pier	with boats	255	1345	
		without boats	1345	255	
quantity	storage tank	one	261	356	
		two	498	119	
		three	552	65	
		four	540	77	
	wast. tr. plant	one	78	724	
		two	758	44	
		three	792	10	
		four	778	24	
	basketball court	one	383	340	
		two	437	286	
		three	702	21	
		half	662	61	
		two-halves	708	15	
	shape	swimming pool	rectangular	299	261
			oval	508	52
			kidney-shaped	313	247
river	curved	623	177		
	straight	177	623		
road	cross	800	800		
	round	800	800		

1177 **Table 5** reports the statistics of xView2-CIR. This dataset should be viewed as a first evaluation benchmark for change-centric RSCIR rather than a complete operational disaster-monitoring dataset. Some disaster categories are small, which motivates reporting both macro-averaged and overall

**Table 5**

**xView2-CIR statistics.** Number of composed queries per disaster type and associated textual modifier. #Queries denotes the number of pre-event composed queries. Each query has exactly one positive match (the post-event image of the *same location* under the *target disaster*), hence #Positives = 1 per query.

DISASTER	TEXT QUERY	#QUERIES	#POSITIVES
hurricane	post-hurricane	147	1
wildfire	post-fire	98	1
flood	post-flood	28	1
tsunami	post-tsunami	9	1
earthquake	post-earthquake	5	1
volcano	post-volcano	4	1

1182 metrics and cautions against over-interpreting fine-grained differences on rare categories.

1184 **A.2. Vision–language backbones**

1185 We evaluate six vision–language models, all using a ViT-L/14 visual backbone:

- 1186 • **CLIP LAION-2B**: A CLIP model trained on LAION-2B [73], a dataset of 2.3 billion web-collected image–text pairs. We use the publicly released laion2b\_s32b\_b82k checkpoint from OpenCLIP [16].
- 1187 • **RemoteCLIP** [51]: A remote-sensing-adapted CLIP model initialized from OpenAI CLIP [67] and fine-tuned on image–text pairs derived from annotated remote sensing datasets using synthetic captions.
- 1188 • **OpenAI CLIP** [67]: The original CLIP model released by OpenAI, trained on 400M web image–text pairs.
- 1189 • **SigLIP** [105]: A vision–language model trained with a sigmoid-based contrastive loss instead of the standard softmax contrastive loss. We use the SigLIP model trained on WebLI.
- 1190 • **CLIP LAION-RS** [94]: A remote-sensing-adapted CLIP model initialized from OpenAI CLIP and fine-tuned on LAION-RS, a 726K-image remote sensing subset of LAION-2B.
- 1191 • **SkyCLIP-50** [94]: A remote-sensing-adapted CLIP model initialized from OpenAI CLIP and fine-tuned on SkyScript-50. SkyScript-50 contains 2.6M image–text pairs constructed by linking geo-referenced satellite imagery from Google Earth Engine [24] with OpenStreetMap [27] annotations.

1212 **A.3. Vocabulary construction**

1213 Some evaluated methods rely on a vocabulary or textual memory. To better align these methods with EO semantics, we generate a family of remote-sensing-specific vocabularies, denoted as RSText, using the prompt shown below. We use vocabularies of different sizes in the analysis and construct HybridText-23k by merging RSText-2k with Open Images v7.

Vocabulary generation prompt

I need a fine-grained, diverse vocabulary list suitable for detailed land-use/land-cover (LULC), remote sensing, and object-detection datasets. The vocabulary should be explicitly descriptive, detailed, and balanced across multiple thematic categories, including but not limited to:

- Urban infrastructure
- Cultural and historical sites
- Recreational and tourism areas
- Transportation and logistics
- Construction and housing
- Education and healthcare
- Natural features and ecosystems
- Biodiversity and wildlife habitats
- Sustainability initiatives and eco-friendly technologies
- Agriculture, farming, and traditional land-use practices
- Marine and aquatic environments

Each vocabulary entry should ideally be short (1–4 words), clear, explicit, self-contained, and suitable for remote sensing imagery annotation. Provide vocabulary entries systematically in batches of 100, carefully ensuring thematic diversity and granularity until a total of at least 2000 entries is reached.

1220

**A.4. Additional details for adapted methods**

This section provides additional implementation details for the adapted composed image retrieval methods evaluated in the main paper. Unless otherwise stated, we follow the official implementations or default settings from the original papers and introduce only the minimal adaptations needed for EO imagery.

*CompoDiff* is a compositional diffusion framework [69] that samples an image-conditioned text embedding via a denoising diffusion process in the joint CLIP embedding space. Sampling is conditioned on (i) the query image embedding, (ii) the composed text prompt (e.g., red {\*}), and (iii) a negative prompt (e.g., low quality). We evaluate CompoDiff in its native configuration and, where applicable, under the hybrid protocol described above. To better preserve identity/class information, we blend the sampled embedding with the original image embedding using a source weight parameter. We report results with 10 diffusion steps and source weight 0.4.

1229

*Pic2Word* learns a lightweight MLP that maps image embeddings to a textual embedding suitable for composition, enabling image-to-text inversion with a single forward pass. For retrieval, we use the template {attribute\_value} {\*}, where \* denotes the inverted image token. We evaluate *Pic2Word* in its standard setting and, when relevant, under the hybrid protocol described above.

*SEARLE* performs test-time textual inversion by optimizing a pseudo token such that its text-encoder embedding matches the query image embedding, regularized by an LLM-guided loss. To ground *SEARLE* in RS, we replace the Open Images v7 [40] vocabulary with our domain-specific HybridText-23k vocabulary and use the RS-adapted prompt template a satellite image of {concept} that. For retrieval, we use the composed text template {attribute\_value} {\$}, where the attribute value (e.g., red) modifies the learned placeholder token \$. We report results with 200 optimization steps, learning rate 0.002, and GPT-loss weight  $\lambda_{\text{GPT}}=0.25$ .

*CIReVL* is a caption-guided composition pipeline that converts the query image into text and then edits it to reflect the modifier. It consists of: (i) captioning the query image, (ii) editing the caption with an instruction-tuned LLM given the textual modifier, and (iii) encoding the edited description with the retrieval backbone for text-to-image search. We use BLIP-2 [43] (blip2-opt-2.7b) as the captioner and LLaMA-2 7B [86] (Llama-2-7b-chat-hf) as the editor. To ground *CIReVL* to PatternCom, we design a domain-specific modifier prompt covering the six attribute types to enforce faithful attribute changes while preserving scene semantics.

1268

CIReVL modifier prompt for PatternCom

You are given a high-resolution satellite image caption describing the “Image Content”, along with a single-word “Instruction” that specifies a modification to apply to the scene. Generate a complete, natural-language “Edited Description” that integrates the modification while preserving all other aspects of the original content. The instruction will always belong to one of six attribute types: color, context, density, existence, quantity, shape. Example (shape):

- Image Content: a satellite image of a rectangular swimming pool in a resort.
- Instruction: oval
- Edited Description: a satellite image of an oval swimming pool in a resort.

*MagicLens* augments an OpenAI CLIP backbone with a transformer based compositional head trained for image–text composition under supervised triplets. We use the official pretrained model (ViT-L/14) and apply it directly to both queries and database images, without additional adaptation.

1274

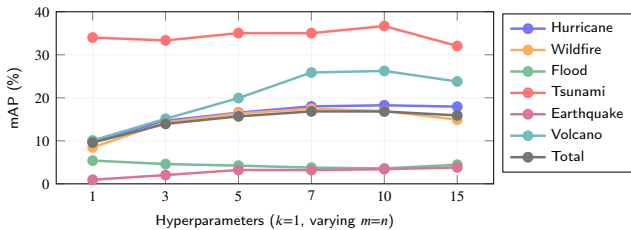
1275 Retrieval is performed using the joint embeddings produced 1318  
 1276 by the MagicLens head. 1319

1277 *WeiCom* is a training-free score-fusion approach that combines 1320  
 1278 image-to-image and text-to-image similarities after 1321  
 1279 calibrating them to a comparable scale. We report results 1322  
 1280 with  $\lambda=0.5$ , corresponding to equal modality contribution. 1323

1281 *BASIC* is a training-free composition method that 1324  
 1282 combines image-to-image and text-to-image similarities 1325  
 1283 after centering and semantic projection, optionally applying 1326  
 1284 image-side query expansion, and fusing modalities with a 1327  
 1285 Harris-regularized multiplicative score. To ground BASIC 1328  
 1286 in RS, we augment  $C_+$  with RSText-150. We report results 1329  
 1287 using 250 principal components for the semantic projection, 1330  
 1288 contrastive scaling  $\alpha=0.2$ , query expansion with 25 nearest 1331  
 1289 neighbors, and Harris regularization weight  $\lambda_{\text{Harris}}=0.1$ .

1290 *FreeDom* is a memory-based textual inversion method that 1332  
 1291 constructs an interpretable text representation of the query 1333  
 1292 image using a large vocabulary and proxy-image expansion. 1334  
 1293 Given a query image and text, it first retrieves  $k$  proxy images 1335  
 1294 via image similarity, then retrieves the top- $n$  vocabulary 1336  
 1295 labels per proxy, and retains the most frequent  $m$  labels 1337  
 1296 overall. Each retained label  $w$  is composed with the modifier 1338  
 1297 using  $\{\text{attribute\_value}\} \{w\}$ , and the resulting text embed- 1339  
 1298 dings are fused by frequency-weighted averaging to form the 1340  
 1299 final query representation. We ground the textual memory 1341  
 1300 to remote sensing using HybridText-23k. We report results 1342  
 1301 using  $k=20$ ,  $n=7$ , and  $m=7$ . 1343

1302 The evaluated methods span a broad deployment spec- 1344  
 1303 trum: some are fully training-free and lightweight at infer- 1345  
 1304 ence (e.g., WeiCom, BASIC), some require large auxiliary 1346  
 1305 vocabularies or proxy retrieval (e.g., FreeDom), some rely 1347  
 1306 on captioning and LLM editing (e.g., CIREVL), and others 1348  
 1307 incur heavier per-query optimization or sampling cost (e.g., 1349  
 1308 SEARLE, CompoDiff). This diversity is important when 1350  
 1309 assessing their suitability for operational EO archives. 1351



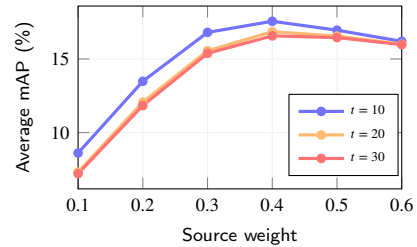
1310 **Figure 11: Impact of frequency-based textual expansion on** 1344  
 1311 **FreeDom with OpenAI CLIP on xView2-CIR with  $k=1$  and** 1345  
 1312 **increasing  $m=n$ . Aggregating multiple textual labels improves** 1346  
 1313 **mAP across most disaster types, especially hurricane and** 1347  
 1314 **wildfire. Performance peaks at  $m=n=7$ , after which further** 1348  
 1315 **expansion yields diminishing returns.** 1349

1313 **A.5. Additional sensitivity and ablation analyses** 1348

1314 *Textual expansion on xView2-CIR.* Figure 11 comple- 1349  
 1315 ments the main-paper analysis of FreeDom textual expan- 1350  
 1316 sion by reporting results on xView2-CIR. Visual expansion 1351  
 1317 is disabled ( $k=1$ ), while the number of aggregated textual

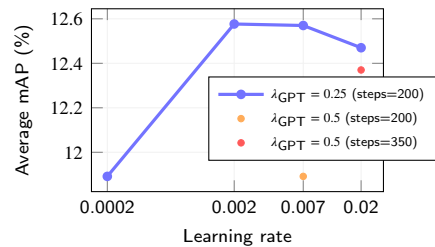
concepts is varied as  $m=n$ . Moderate textual expansion 1318  
 improves performance, with the best Total score obtained 1319  
 at  $m=n=7$ . Larger expansions provide diminishing returns, 1320  
 suggesting that overly long textual surrogates can introduce 1321  
 noisy or off-target concepts. 1322

*CompoDiff hyperparameter sensitivity.* Figure 12 ana- 1323  
 lyzes CompoDiff sensitivity to the number of diffusion 1324  
 timesteps  $t$  and the source weight used to mix the refer- 1325  
 ence image with the generated proxy. Performance improves 1326  
 when moving from very small source weights to the mid- 1327  
 range (0.3–0.5), indicating that both the original image and 1328  
 the generative signal are useful for composition. Moderate 1329  
 diffusion budgets ( $t=10-20$ ) are generally sufficient, while 1330  
 larger budgets provide limited additional gains. 1331



1332 **Figure 12: CompoDiff hyperparameter sweep with SkyCLIP-** 1333  
 1334 **50 on PatternCom. We report average mAP as a function of** 1335  
 1336 **source weight for different numbers of diffusion timesteps.** 1337

*SEARLE hyperparameter sensitivity.* Figure 13 evalu- 1338  
 ates SEARLE under different learning rates, inversion steps, 1339  
 and GPT mixing coefficients  $\lambda_{\text{GPT}}$ . Performance is relatively 1340  
 stable around  $\lambda_{\text{GPT}}=0.25$  and moderate learning rates, while 1341  
 very small learning rates degrade results. Larger GPT mix- 1342  
 ing does not provide consistent gains. Overall, SEARLE is 1343  
 moderately sensitive to optimization settings but remains 1344  
 stable around prior-work defaults. 1345



1342 **Figure 13: SEARLE hyperparameter sweep with CLIP LAION-** 1343  
 1344 **2B on PatternCom. We report average mAP as a function of** 1345  
 1346 **learning rate, number of inversion steps, and GPT mixing** 1347  
 1348 **coefficient  $\lambda_{\text{GPT}}$ .** 1349

*Ablation of BASIC components* Table 6 and 7 analyze 1350  
 the contribution of BASIC components on PatternCom and 1351  
 xView2-CIR. Across both datasets, removing *centering* or 1352  
*semantic projection* causes the largest drops, showing that 1353  
 modality calibration and semantic subspace alignment are 1354  
 the main drivers of BASIC’s performance. Other compo- 1355  
 nents, including Harris weighting, min-based normalization, 1356  
 text contextualization, and query expansion, provide smaller 1357

**Table 6**

**BASIC component ablation on PatternCom with OpenAI CLIP.** We report attribute-wise and average mAP (%) when removing components from the full pipeline: mean centering (**Cen.**), min-based normalization (**Norm.**), Harris weighting (**Har.**), text contextualization (**Con.**), semantic projection (**Proj.**), and query expansion (**Q. Ex.**).

Cen.	Norm.	Har.	Con.	Proj.	Q. Ex.	Color	Context	Density	Existence	Quantity	Shape	Avg.
✓	✓	✓	✓	✓	✓	38.31	24.53	26.35	30.94	18.03	35.73	<b>28.98</b>
✓	✓	✗	✓	✓	✓	36.69	22.87	24.34	28.40	16.14	32.74	26.86
✓	✗	✗	✓	✓	✓	26.57	30.51	24.63	11.93	26.08	23.27	23.83
✓	✓	✓	✗	✓	✓	34.60	24.56	19.17	12.81	26.44	17.03	22.44
✓	✓	✓	✓	✗	✓	23.22	24.41	18.97	3.11	24.40	20.19	19.05
✗	✓	✓	✓	✓	✓	28.90	18.95	22.77	8.02	12.57	16.33	17.92
✓	✓	✓	✓	✓	✗	37.16	24.51	24.93	28.62	16.57	34.22	27.67

**Table 7**

**BASIC component ablation on xView2-CIR with OpenAI CLIP.** We report disaster-wise mAP (%), macro-average mAP (Avg.), and overall mAP (Total) when removing components from the full pipeline. Component abbreviations follow Table 6.

Cen.	Norm.	Har.	Con.	Proj.	Q. Ex.	Hurricane	Fire	Flood	Tsunami	Earthquake	Volcano	Avg.	Total
✓	✓	✓	✓	✓	✓	9.36	7.74	7.54	19.52	1.18	13.96	9.88	8.88
✓	✓	✓	✓	✓	✗	22.71	18.49	14.59	46.15	1.69	59.20	<b>27.14</b>	<b>21.38</b>
✓	✓	✗	✓	✓	✗	21.16	17.45	14.60	47.63	1.79	51.04	25.61	20.18
✓	✗	✗	✓	✓	✗	16.39	17.08	11.11	29.13	1.14	64.17	23.17	16.90
✓	✓	✗	✗	✓	✗	18.48	16.69	14.38	32.37	2.62	57.25	23.63	18.17
✓	✓	✓	✓	✗	✗	16.03	14.74	13.68	24.25	1.26	59.01	21.49	15.96
✗	✓	✓	✓	✓	✗	0.68	0.71	1.24	17.02	1.39	6.29	4.56	1.34

and more dataset-dependent gains. Notably, query expansion is harmful on xView2-CIR, where relevance requires the same scene/location and visually similar proxies often introduce off-location evidence. Overall, BASIC’s gains primarily come from representation calibration and semantic projection, while auxiliary heuristics play a secondary role.

**A.6. Sensitivity to text query phrasing**

We study how composed retrieval methods respond to alternative textual formulations of the modifier. For PatternCom, we test three templates applied to each attribute value  $v$ : R1 uses a *being* formulation (e.g., being blue, being rectangular-shaped); R2 uses a *having* formulation (e.g., having blue color, having rectangular shape); and R3 uses longer, relational descriptions such as with the main object modified to have blue color or with the main object modified to have rectangular shape. As shown in Table 8, the original plain attribute word (e.g., blue, water, dense) is generally the most reliable choice, especially for stronger multimodal methods. More verbose templates tend to reduce performance, suggesting that additional function words can dilute the compact attribute signal.

For xView2-CIR, the plain query is post- $\{\text{disaster}\}$ , while R1 uses impact-oriented rephrasing (e.g., burned area, flooded area) and R2 uses disaster-explicit formulations (e.g., fire-affected region, seismic damage). As shown in Table 9, prompt sensitivity is more method-dependent in this change-centric setting. More descriptive formulations can benefit methods that rely more strongly on semantic text structure, particularly BASIC, while other methods are more sensitive to distribution shifts introduced by rewording. Overall, these results indicate that prompt design is

**Table 8**

**Text query rephrasing on PatternCom with CLIP LAION-2B.** We report average mAP (%) for the original modifier (Plain) and three rephrased variants (R1–R3).

METHOD	Average			
	Plain	R1	R2	R3
Text-only	5.57	<b>9.16</b>	7.07	3.52
Text+Image	16.41	<b>17.70</b>	17.48	15.10
Text×Image	21.00	<b>19.03</b>	<b>21.58</b>	16.76
WeiCom	21.75	20.86	<b>23.31</b>	16.38
BASIC	16.47	<b>17.50</b>	17.73	17.34
FreeDom	<b>38.59</b>	33.06	31.48	29.53

**Table 9**

**Text query rephrasing on xView2-CIR with OpenAI CLIP.** We report macro-average mAP (Average) and overall mAP (Total) for the original query (Plain) and two rephrased variants (R1–R2).

METHOD	Average			Total		
	Plain	R1	R2	Plain	R1	R2
Text-only	<b>7.57</b>	4.53	4.30	<b>2.97</b>	1.76	2.06
Text+Image	<b>17.73</b>	15.27	15.63	<b>12.14</b>	9.14	10.22
Text×Image	<b>25.94</b>	14.97	12.13	<b>16.35</b>	10.38	9.65
WeiCom	<b>30.67</b>	24.39	19.22	<b>21.40</b>	12.58	13.13
BASIC	<b>27.14</b>	17.39	17.36	<b>21.38</b>	18.56	18.17
FreeDom	<b>17.23</b>	9.86	11.98	<b>16.84</b>	10.83	12.19

dataset- and method-dependent: compact modifiers are effective for controlled attribute edits, whereas change-centric retrieval may benefit from more explicit event or impact descriptions.