

Compositional Preference Learning for Composed Person Retrieval

Anonymous CVPR submission

Paper ID *****

Abstract

001 In aerial surveillance, person retrieval is challenging due
 002 to low resolution, large viewpoint changes, and appearance
 003 ambiguity, which limit appearance-only matching. These
 004 conditions motivate multimodal person retrieval settings in
 005 which identity cues from a reference image are combined
 006 with language-based descriptions of appearance change.
 007 We study this problem through composed person retrieval
 008 (CPR), where a target image is retrieved given a reference
 009 image and a modification text. We propose a simple frame-
 010 work that models compositional preference by construct-
 011 ing mismatched compositions via cross-sample replacement
 012 of either images or texts, and training the model to rank
 013 the original composition above these variants. This objec-
 014 tive encourages the model to remain sensitive to changes
 015 in either the modification text or the reference identity, en-
 016 abling finer-grained compositional reasoning. Our method
 017 achieves state-of-the-art performance on the ITCPR bench-
 018 mark, surpassing the previous best supervised CPR base-
 019 line by 2.18% in Rank-1 and 1.80% in mAP. These results
 020 demonstrate that explicitly modeling compositional prefer-
 021 ence is an effective strategy for composed person retrieval
 022 and a promising direction for challenging surveillance sce-
 023 narios, particularly aerial surveillance.

024 1. Introduction

025 In aerial and cross-view surveillance scenarios [4, 9, 14,
 026 17], person retrieval is particularly challenging due to low
 027 resolution, large viewpoint variations, and appearance am-
 028 biguity, which often render appearance-only matching un-
 029 reliable. These challenges motivate retrieval formulations
 030 that integrate visual identity cues with language-based de-
 031 scriptions of appearance changes. In this work, we study
 032 this problem in the framework of composed person retrieval
 033 (CPR), which aims to retrieve a target person image from
 034 a gallery given a multimodal query composed of a refer-
 035 ence image and a modification text. Unlike conventional
 036 person retrieval, CPR [7, 8, 16] is designed for scenarios
 037 where the target person should preserve the core identity of

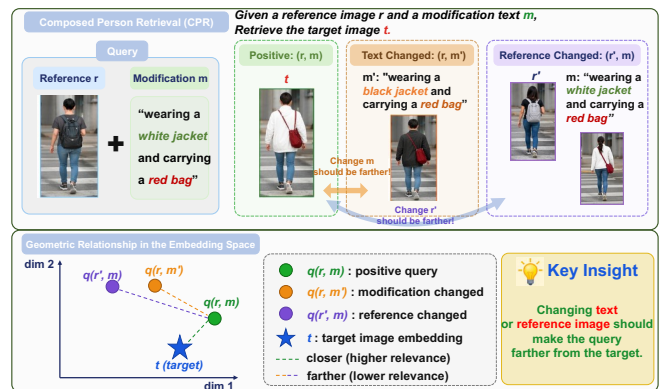


Figure 1. Relational structure in composed person retrieval. The original query should remain more compatible with the target image than mismatched variants constructed by changing either the modification text or the reference image.

the reference while reflecting a specific appearance change described in natural language, such as changes in clothing, accessories, or other visible attributes. As a result, CPR requires both identity preservation and fine-grained compositional reasoning across vision and language.

The central difficulty of CPR lies in distinguishing *which components* of the multimodal query govern identity preservation and *which part* specifies appearance change, rather than merely aligning the reference image, the text, and the target image at a coarse level. *However, existing CPR methods mainly learn from positive triplets, without explicitly exploiting the richer relational structure among triplets within a mini-batch.*

In this paper, we revisit CPR from the perspective of relational compositional supervision. Our key observation is that samples within a training batch naturally allow the construction of informative mismatched compositions: one can replace the modification text of a triplet with that of another, or replace the reference image while keeping the modification text fixed. Although these variants are semantically close to the original query, they should be less compatible with the corresponding target. This provides a simple yet effective way to introduce finer-grained supervision beyond

061 standard triplet matching. As illustrated in Fig. 1, changing
062 either the modification text or the reference identity should
063 move the composed query farther away from the correct tar-
064 get in the embedding space.

065 Based on this intuition, we propose a straightforward
066 framework CPR framework: **CO**mpositional **P**reference
067 learning for composEd person **R**etrieval (COPER). COPER
068 strengthens model training by mining mismatched compo-
069 sitions within each mini-batch and introduces a Composi-
070 tional Preference Loss that encourages the original query-
071 target pair to be ranked above its corrupted counterparts.
072 By doing so, the model is trained not only to match valid
073 compositions, but also to distinguish them from subtly mis-
074 matched alternatives that differ in either appearance modi-
075 fication or identity cue.

076 Our main contributions are summarized as follows:

- 077 • We identify the lack of fine-grained compositional super-
078 vision as a key limitation in existing CPR methods, and
079 revisit the task from a relational reasoning perspective
080 within a mini-batch.
- 081 • We propose a Compositional Preference Loss based on
082 mismatched composition mining, which constructs mis-
083 matched triplets by swapping the reference image or mod-
084 ification text across samples and encourages the model to
085 rank the original composition higher.
- 086 • We validate the effectiveness of our method on the
087 ITCPR benchmark, achieving state-of-the-art perfor-
088 mance among supervised CPR approaches.

089 2. Related Work

090 2.1. Composed Image Retrieval

091 Composed image retrieval (CIR) extends standard retrieval
092 by using a multimodal query composed of a reference im-
093 age and a modification text. Existing CIR methods can be
094 divided into supervised and zero-shot approaches. Super-
095 vised methods learn from triplet annotations, while zero-
096 shot methods avoid such supervision by leveraging pre-
097 trained vision-language models and pseudo-word based
098 composition. Representative zero-shot CIR methods in-
099 clude Pic2Word [10], SEARLE [1], and LinCIR [2]. How-
100 ever, these methods are mainly designed for natural images
101 and do not explicitly address the fine-grained identity sen-
102 sitivity required in person retrieval.

103 2.2. Composed Person Retrieval

104 Composed person retrieval (CPR) extends composed im-
105 age retrieval to person search, where the goal is to retrieve
106 images of the same person after appearance changes de-
107 scribed by text. Word4Per [8] is the first work to formu-
108 late CPR and introduce the ITCPR benchmark. I2ID [16]
109 further improves this line by suppressing non-identity ap-
110 pearance cues that may conflict with the modification text.

FAFA [7] advances CPR by introducing SynCPR, a large-
scale synthetic training dataset, together with an adaptive
feature alignment framework.

114 3. Method

115 3.1. Task Formulation

116 Given a dataset $\mathcal{D} = \{(r_i, m_i, t_i)\}_{i=1}^M$, where r_i denotes a
117 reference image, m_i a modification text, and t_i the corre-
118 sponding target image, the goal is to learn a retrieval model
119 that composes the identity-related visual information from
120 r_i with the semantic modification described by m_i . For a
121 mini-batch of size N sampled from \mathcal{D} , we denote the sets of
122 reference images, modification texts, and target images by
123 $\mathcal{R} = \{r_i\}_{i=1}^N$, $\mathcal{M} = \{m_i\}_{i=1}^N$, and $\mathcal{T} = \{t_i\}_{i=1}^N$, respec-
124 tively. Given a query pair (r_i, m_i) , the model aims to re-
125 trieve the corresponding target image t_i from the candidate
126 set \mathcal{T} such that the retrieved image reflects the modification
127 in m_i while preserving the identity-related characteristics
128 of r_i .

129 3.2. Overall framework

130 Our COPER framework is built on the Q-Former architec-
131 ture from BLIP-2 [5], which enables efficient multimodal
132 interaction between image and text representations while
133 requiring only lightweight trainable modules. Let $E_{\text{img}}(\cdot)$
134 denote the image encoder and $Q(\cdot)$ the Q-Former. For each
135 triplet (r_i, m_i, t_i) , the reference image r_i and modification
136 text m_i are jointly processed to produce a composed query
137 embedding, while the target image t_i is encoded into a tar-
138 get representation:

$$139 q_i = Q(E_{\text{img}}(r_i), m_i), \quad v_i = Q(E_{\text{img}}(t_i)). \quad (1)$$

140 Here, q_i denotes the final text-side [CLS] embedding pro-
141 duced from the joint encoding of r_i and m_i , which serves
142 as the composed query representation. In contrast, $v_i =$
143 $\{v_i(l)\}_{l=1}^{N_q}$ denotes the set of N_q target-side query tokens
144 output by the Q-Former for t_i . Thus, q_i captures both the
145 identity cues of r_i and the appearance modification spec-
146 ified by m_i , while v_i provides fine-grained target-side to-
147 ken features for relevance scoring. Given the query pair
148 (r_i, m_i) , the model is trained so that q_i is highly compati-
149 ble with the corresponding target token set v_i while remain-
150 ing distinguishable from other candidate targets in the mini-
151 batch.

152 Following FAFA [7], we use the FDA loss (\mathcal{L}_{FDA}) as the
153 base retrieval objective. In addition, we introduce a Composi-
154 tional Preference Loss to provide additional supervision
155 on the relational structure among samples within a mini-
156 batch. The detailed formulation of this loss is described in
157 the next subsection.

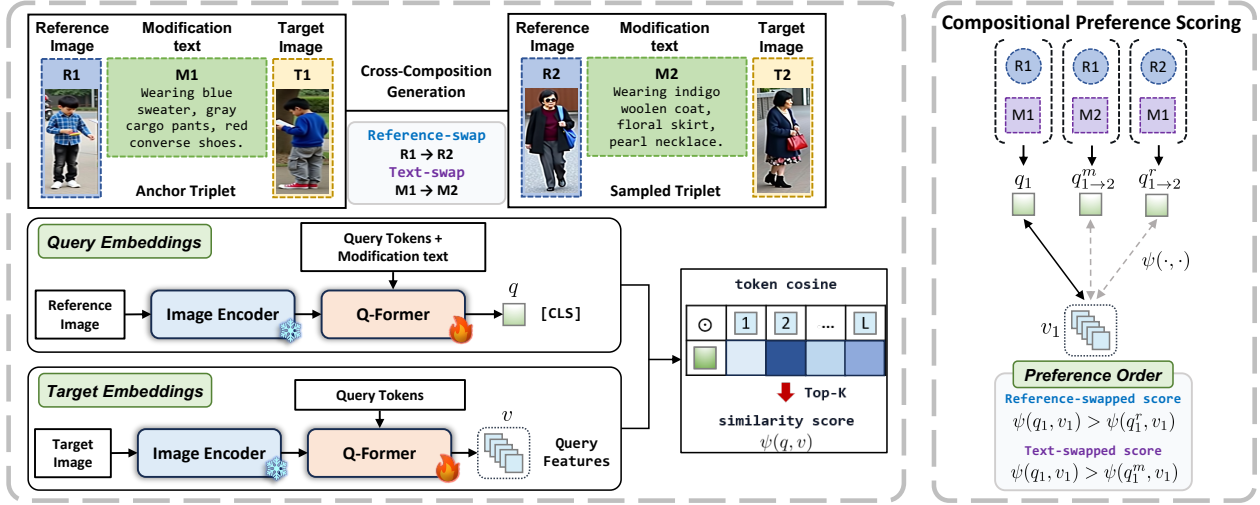


Figure 2. Overview of the COPER framework. COPER encodes the reference image, modification text, and target image with Q-Former, constructs cross-compositions within each mini-batch, and is trained with fine-grained alignment and Compositional Preference supervision.

3.3. Compositional Preference Loss

While the FDA loss encourages fine-grained alignment between a composed query and its target image, it does not explicitly enforce sensitivity to which component of the query has changed. To provide such additional supervision, we introduce a Compositional Preference Loss based on mismatched composition mining.

First, for each sample i in a mini-batch, we randomly select another sample $j \neq i$ from the same batch and construct two cross-composed queries by replacing either the modification text or the reference image:

$$q_{i \rightarrow j}^m = Q(E_{\text{img}}(r_i), m_j), \quad (2)$$

$$q_{i \rightarrow j}^r = Q(E_{\text{img}}(r_j), m_i). \quad (3)$$

Here, $q_{i \rightarrow j}^m$ keeps the reference image fixed while replacing the modification text, whereas $q_{i \rightarrow j}^r$ keeps the modification text fixed while replaces the reference image. These cross-composed queries remain semantically close to the original composition, but they should be less compatible with the corresponding target image.

Next, we measure the compatibility between a query embedding q and a target representation v using token-level cosine similarities followed by top- K aggregation:

$$\psi(q, v) = \frac{1}{K} \sum_{k=1}^K \text{TopK}_k \left(\left\{ \frac{q^\top v_i(l)}{\|q\| \|v_i(l)\|} \right\}_{l=1}^{N_q} \right), \quad (4)$$

where $v(l)$ denotes the l -th target token and $\text{TopK}_k(\cdot)$ denotes the k -th largest value. Based on this measure, the compatibility score of the original composition is defined as

$$s_i^+ = \psi(q_i, v_i). \quad (5)$$

We then compute the compatibility scores of the two cross-composed queries with respect to the corresponding target:

$$s_{i \rightarrow j}^{x,i} = \psi(q_{i \rightarrow j}^x, v_i), \quad x \in \{m, r\}. \quad (6)$$

Finally, we encourage the original composition to be more compatible with the target image than either cross-composed variants:

$$\mathcal{L}_{\text{pref}}^x = -\log \sigma \left(\frac{s_i^+ - s_{i \rightarrow j}^{x,i}}{\tau} \right), \quad x \in \{m, r\}, \quad (7)$$

where $\sigma(\cdot)$ denotes the sigmoid function and τ is a temperature parameter. The final Compositional Preference Loss is defined as

$$\mathcal{L}_{\text{CP}} = \frac{1}{N} \sum_{i=1}^N (\mathcal{L}_{\text{pref}}^m + \mathcal{L}_{\text{pref}}^r). \quad (8)$$

This objective encourages the model to remain sensitive to changes in either the modification text or the reference identity, enabling finer-grained compositional reasoning. The final training objective is

$$\mathcal{L} = \mathcal{L}_{\text{FDA}} + \lambda \mathcal{L}_{\text{CP}}, \quad (9)$$

where \mathcal{L}_{FDA} denotes the base FDA loss and λ controls the contribution of the proposed compositional preference loss.

4. Experiments

4.1. Datasets

We follow the experimental protocol of FAFA [7] and use separate datasets for training and evaluation. For training,

Table 1. Comparison of our method with representative approaches from related retrieval settings, including ZSCIR, CIR, ZSCPR, and CPR. “Combination” indicates that both the reference image and the modification text are used as the query. “PR” denotes the journal *Pattern Recognition*. The best overall results are highlighted in bold.

| Domain | Method | Venue | Pretraining Data | Setting | Rank-1 | Rank-5 | Rank-10 | mAP |
|--------|--------------------|------------|------------------|-------------|--------------|--------------|--------------|--------------|
| ZSCIR | Pic2Word [10] | CVPR’23 | CC3M [11] | Combination | 21.21 | 37.15 | 44.51 | 29.11 |
| | CoVR-BLIP [13] | AAAI’24 | WebVid-CoVR [13] | Combination | 26.75 | 47.68 | 56.36 | 36.49 |
| | LinCIR (ViT-G) [2] | CVPR’24 | – | Combination | 23.93 | 44.46 | 53.18 | 33.95 |
| CIR | CaLa [3] | SIGIR’24 | SynCPR [7] | Combination | 39.33 | 60.85 | 68.66 | 49.29 |
| | SPRC [15] | ICLR’24 | SynCPR [7] | Combination | 42.27 | 61.81 | 69.35 | 51.62 |
| ZSCPR | Word4Per [8] | ARXIV | CUHK-PEDES [6] | Combination | 44.23 | 64.76 | 72.17 | 54.44 |
| | I2ID [16] | PR’26 | CUHK-PEDES [6] | Combination | 46.41 | 66.49 | 73.52 | 55.53 |
| CPR | FAFA [7] | NeurIPS’25 | SynCPR [7] | Combination | 45.46 | 66.44 | 73.30 | 55.15 |
| | COPER(Ours) | - | SynCPR [7] | Combination | 47.64 | 67.85 | 74.98 | 56.95 |

208 we use the synthetic SynCPR dataset [7]. Each sample
 209 is constructed via an automatic multimodal data synthesis
 210 pipeline and consists of a reference image, a relative cap-
 211 tion describing the modification, and the corresponding tar-
 212 get image. SynCPR provides large-scale supervision with
 213 diverse identities, clothing variations, poses, and scenes,
 214 thereby enabling fine-grained alignment between visual and
 215 textual cues.

216 For evaluation, we use the ITCPR benchmark from
 217 Word4Per [8], a manually annotated test set introduced for
 218 CPR. ITCPR is designed to reflect realistic retrieval sce-
 219 narios where both a reference image and a relative textual
 220 description are available at query time.

221 4.2. Implementation Details

222 We train COPER on SynCPR for 10 epochs using AdamW,
 223 with an initial learning rate of $4e-6$ and a global batch size
 224 of 256. All experiments are conducted on a single NVIDIA
 225 B200 GPU. The model adopts EVA-CLIP-g [12] as the im-
 226 age encoder with an input resolution of 224×224 , along
 227 with a Q-Former containing 32 learnable query tokens. For
 228 optimization, we set the weight of the proposed Composi-
 229 tional Preference Loss to $\lambda = 1.0$ and the temperature
 230 parameter to $\tau = 0.07$.

231 4.3. Comparison with State-of-the-Art

232 We compare our method with representative approaches
 233 from related retrieval settings, including ZSCIR, CIR,
 234 ZSCPR, and CPR, as summarized in Table 1. General com-
 235 posed retrieval methods, including both zero-shot and su-
 236 pervised CIR approaches, remain clearly inferior to CPR-
 237 oriented methods, suggesting that person retrieval requires
 238 stronger identity-aware compositional reasoning. At the
 239 same time, zero-shot CPR methods such as Word4Per
 240 and I2ID show competitive performance, highlighting the

Table 2. Ablation study on the effect of the proposed Compositional Preference loss \mathcal{L}_{CP} .

| Method | Rank-1 | Rank-5 | Rank-10 | mAP |
|------------------------|---------------|---------------|---------------|---------------|
| w/o \mathcal{L}_{CP} | 46.00 | 65.99 | 73.12 | 55.32 |
| w/ \mathcal{L}_{CP} | 47.64 (+1.64) | 67.85 (+1.86) | 74.98 (+1.86) | 56.95 (+1.63) |

strong transferability of pretrained vision-language mod- 241
 els to this task. Our method achieves the best overall per- 242
 formance, surpassing FAFA by +2.18% points in Rank-1, 243
 +1.41% points in Rank-5, +1.68% points in Rank-10, and 244
 +1.80% points in mAP, demonstrating the effectiveness of 245
 our approach for CPR. 246

247 4.4. Ablation Study

To verify the effectiveness of the proposed Compositional 248
 Preference Loss, we compare the full model with a vari- 249
 ant trained without \mathcal{L}_{CP} on ITCPR. As shown in Table 2, 250
 removing \mathcal{L}_{CP} leads to consistent drops across all metrics, 251
 confirming its role in learning identity-aware compositional 252
 relations. 253

254 5. Conclusion

In this paper, we address composed person retrieval by 255
 explicitly modeling its identity-sensitive and composition- 256
 aware nature. We propose a straightforward framework 257
 based on mismatched composition mining and a Composi- 258
 tional Preference Loss, which improves the discrimination 259
 between valid and corrupted compositions. We believe this 260
 framework can serve as a useful step toward multimodal 261
 person retrieval in challenging surveillance scenarios, in- 262
 cluding aerial and cross-view settings. 263

264

References

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

- [1] Alberto Baldrati, Lorenzo Agnolucci, Marco Bertini, and Alberto Del Bimbo. Zero-shot composed image retrieval with textual inversion. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 15338–15347, 2023. 2
- [2] Geonmo Gu, Sanghyuk Chun, Wonjae Kim, Yoohoon Kang, and Sangdoon Yun. Language-only training of zero-shot composed image retrieval. In *CVPR*, pages 13225–13234, 2024. 2, 4
- [3] Xintong Jiang, Yaxiong Wang, Mengjian Li, Yujiao Wu, Bingwen Hu, and Xueming Qian. Cala: Complementary association learning for augmenting composed image retrieval. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2177–2187, 2024. 4
- [4] Khadija Khaldi, Vuong D. Nguyen, Pranav Mantini, and Shishir Shah. Unsupervised person re-identification in aerial imagery. In *IEEE/CVF Winter Conference on Applications of Computer Vision Workshops (WACVW)*, pages 260–269, 2024. 1
- [5] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 2
- [6] Shuang Li, Tong Xiao, Hongsheng Li, Bolei Zhou, Dayu Yue, and Xiaogang Wang. Person search with natural language description. In *CVPR*, pages 1970–1979, 2017. 4
- [7] Delong Liu, Haiwen Li, Zhaohui Hou, Zhicheng Zhao, Fei Su, and Yuan Dong. Automatic synthetic data and fine-grained adaptive feature alignment for composed person retrieval. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*. 1, 2, 3, 4
- [8] Delong Liu, Haiwen Li, Zhaohui Hou, Zhicheng Zhao, Fei Su, and Yuan Dong. Automatic synthetic data and fine-grained adaptive feature alignment for composed person retrieval. *arXiv preprint arXiv:2311.16515*, 2023. 1, 2, 4
- [9] Huy Nguyen, Kien Nguyen, Sridha Sridharan, and Clinton Fookes. Aerial-ground person re-id. In *IEEE International Conference on Multimedia and Expo (ICME)*, pages 2585–2590, 2023. 1
- [10] Kuniaki Saito, Kihyuk Sohn, Xiang Zhang, Chun-Liang Li, Chen-Yu Lee, Kate Saenko, and Tomas Pfister. Pic2word: Mapping pictures to words for zero-shot composed image retrieval. In *CVPR*, pages 19305–19314, 2023. 2, 4
- [11] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, 2018. 4
- [12] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training techniques for clip at scale. *arXiv preprint arXiv:2303.15389*, 2023. 4
- [13] Lucas Ventura, Antoine Yang, Cordelia Schmid, and G’ul Varol. Covr: Learning composed video retrieval from web video captions. In *AAAI*, pages 5270–5279, 2024. 4
- [14] Lei Wang, Quan Zhang, Junyang Qiu, and Jianhuang Lai. Rotation exploration transformer for aerial person re-identification. In *IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6, 2024. 1
- [15] Xinxing Xu, Yong Liu, Salman Khan, Fahad Khan, Wangmeng Zuo, Rick Siow Mong Goh, Chun-Mei Feng, et al. Sentence-level prompts benefit composed image retrieval. In *ICLR*, 2024. 4
- [16] Guo Yu, Di Wang, Chengwei Yan, Feng Yan, Nan Luo, Yifeng Wang, and Quan Wang. I2id: Disentangling identity features via synchronized masking for zero-shot composed person retrieval. *Pattern Recognition*, 179:113654, 2026. 1, 2, 4
- [17] Shizhou Zhang, Qi Zhang, Yifei Yang, Xing Wei, Peng Wang, Bingliang Jiao, and Yanning Zhang. Person re-identification in aerial imagery. *IEEE Transactions on Multimedia*, 23:281–291, 2021. 1