

Collaborative Group: Composed Image Retrieval via Consensus Learning from Noisy Annotations

Xu Zhang^a, Zhedong Zheng^{b,*}, Linchao Zhu^a, Yi Yang^a

^a*College of Computer Science and Technology, Zhejiang University, Hangzhou, 310058, China*

^b*Faculty of Science and Technology, and Institute of Collaborative Innovation, University of Macau, Macau, 999078, China*

Abstract

Composed image retrieval extends content-based image retrieval systems by enabling users to search using reference images and captions that describe their intention. Despite great progress in developing image-text compositors to extract discriminative visual-linguistic features, we identify a hitherto overlooked issue, triplet ambiguity, which impedes robust feature extraction. Triplet ambiguity refers to a type of semantic ambiguity that arises between the reference image, the relative caption, and the target image. It is mainly due to the limited representation of the annotated text, resulting in many noisy triplets where multiple visually dissimilar candidate images can be matched to an identical reference pair (i.e., a reference image + a relative caption). To address this challenge, we propose the Consensus Network (Css-Net), inspired by the psychological concept that groups outperform individuals. Css-Net comprises two core components: (1) a consensus module with four diverse compositors, each generating distinct image-text embeddings, fostering complementary feature extraction and mitigating dependence on any single, potentially biased compositor; (2) a Kullback-Leibler divergence loss that encourages learning of inter-compositor interactions to promote consensual outputs. During evaluation, the decisions of the four compositors are combined through a weighting scheme, enhancing overall agreement. On benchmark datasets, particularly FashionIQ, Css-Net demonstrates marked improvements. Notably, it achieves significant recall gains, with a 2.77% in-

*Corresponding author.

Email address: zhedongzheng@um.edu.mo (Zhedong Zheng)

crease in R@10 and 6.67% boost in R@50, underscoring its competitiveness in addressing the fundamental limitations of existing methods.

Keywords: Noisy Annotation, Data Ambiguity, Compositional Image Retrieval, Image Retrieval with Text Feedback, Multi-modal Retrieval

1. Introduction

Image retrieval plays a pivotal role in computer vision and proves to be valuable in many applications, such as product search (Guo et al., 2019; Sharma and Vishvakarma, 2019; Guo et al., 2018), internet search (Noh et al., 2017) and fashion retrieval (Liu et al., 2016; Liao et al., 2018). Prevalent image retrieval approaches include image-to-image retrieval (Deng et al., 2019; Fan et al., 2019; Sheng et al., 2020; Hafner et al., 2022) and text-to-image retrieval (Zhen et al., 2019; Zheng et al., 2020; Guerrero et al., 2021; Wang et al., 2022), which endeavor to locate the image of interest using a single image or descriptive texts as a query. Despite significant progress, users often lack a precise search target in advance but instead seek categories, such as shoes or clothing. Therefore, an interactive system is highly desirable to assist users to reconsider their intentions, as depicted in Fig. 1. Hence, Composed image retrieval, which aims to search the image of interest given the composed query consisting of a reference image and a relative caption describing the modification, has attracted great attention (Vo et al., 2019; Chen et al., 2020; Lee et al., 2021; Kim et al., 2021; Wen et al., 2021).

Recent studies addressing the task of composed image retrieval primarily concentrate on extracting discriminative representations from image-text-image triplets. For example, TIRG (Vo et al., 2019), VAL (Chen et al., 2020), and CoSMo (Lee et al., 2021) propose different ways to modify the visual features of the reference image conditioned on the relative caption. TIRG uses a simple gating and residual module, VAL devises a visual-linguistic attention learning framework, and CoSMo introduces the content and style modulators. Additionally, CLVC-Net (Wen et al., 2021) and CLIP4cir (Baldrati et al., 2022) devise more intricate multi-modal fusion modules to accentuate the modifications of the reference image. CLVC-Net uses local-wise and global-wise compositors, while CLIP4cir finetunes the CLIP (Radford et al., 2021) text encoder and trains a combiner network to fuse features.

Despite the significant success, these works fail to address an inherent problem of the composed image retrieval task: the ambiguity of the training

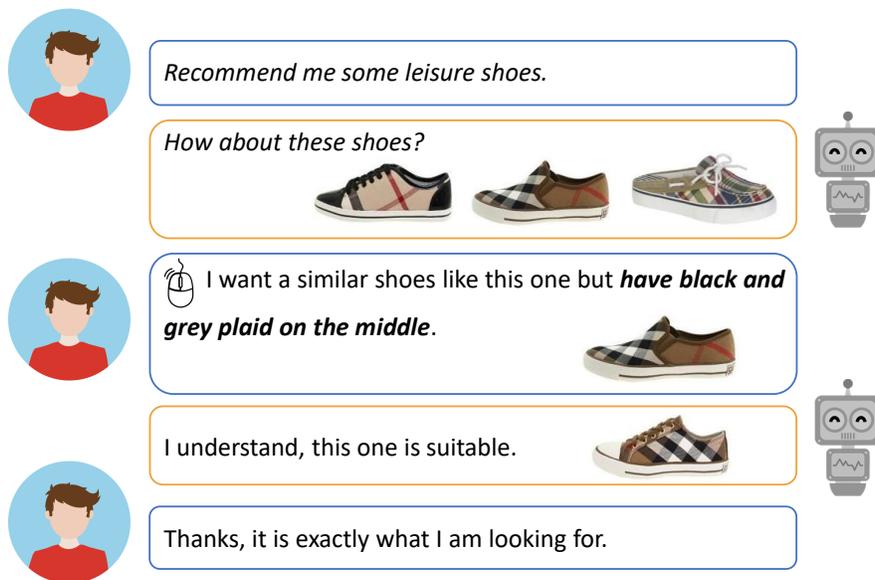


Figure 1: Schematic illustration of the composed image retrieval system. Through using a reference image and a relative caption, the system endeavors to precisely retrieve the intended target image from all candidate images.

32 data triplets, *i.e.*, **triplet ambiguity**. Triplet ambiguity originates from
 33 the annotation process where annotators focus on single data triplet, and
 34 frequently describe simple properties such as color and size, while neglect-
 35 ing more fine-grained details, such as location and style. Consequently, many
 36 noisy triplets exist where candidate images meet the requirement of the com-
 37 posed query but are not annotated as the desired ground-truth target image,
 38 especially when the relative caption is brief. Similar annotation ambiguity
 39 is also observed in pair-wise data (Wray et al., 2021; Falcon et al., 2022)
 40 and remains challenging. As shown in Fig. 2, existing methods treat com-
 41 posed image retrieval as an instance-level retrieval, that is, given a refer-
 42 ence pair (comprising a reference image and a relative caption), only the anno-
 43 tated target image is considered as the correct image to retrieve. In fact,
 44 due to the limitation of the text description, many candidate images within
 45 the dataset are semantically similar to the point of being identical, but are
 46 treated as the negative counterparts, thus producing many noisy triplets.
 47 These noisy triplets compromises the representation learning of the single
 48 compositor, since the metric learning objective in this task aims to push
 49 away these false-negative samples from the composed query. Empirically, we



Figure 2: Illustration of the triplet ambiguity problem. Triplet ambiguity denotes multiple false-negative samples in the dataset as the annotator usually see one triplet with true match (✓) at a time, while neglecting other candidates (?).

50 verify the existence of triplet ambiguity in Sec. 4.2.

51 To relieve the triplet ambiguity problem, we propose a straightforward
 52 and effective Consensus Network (Css-Net) for composed image retrieval, as
 53 illustrated in Fig. 3(a). The key idea underpinning our method to allevi-
 54 ate the triplet ambiguity is “two heads are better than one” in short. To
 55 be more specific, an individual often errs due to the biases caused by noisy
 56 triplets, but groups are less susceptible to making similar mistakes, thereby
 57 circumventing sub-optimal solutions. This is known as the psychological find-
 58 ing (Hinsz, 1990) that groups perform better than individuals on the memory
 59 task. Consequently, we aim to (1) develop a consensus module (group) com-
 60 posed of composers possessing diverse knowledge to jointly make decisions
 61 during evaluation and (2) encourage learning among different composers to
 62 minimize their biases learned on noisy triplets by employing an additional
 63 Kullback Leibler divergence loss (KL loss) (Kullback and Leibler, 1951).

64 Css-Net ensures that the composers possess distinct knowledge in two
 65 ways: • Motivated by the finding (Lin et al., 2017; Miech et al., 2021) that
 66 the image features of high-resolution are semantically weak, while the image
 67 features of low-resolution are semantically strong, we employ two image-text
 68 composers at different depths of the same image encoder, (*i.e.*, block3 and
 69 block4 of the ResNet (He et al., 2016)). The former focuses more on detailed
 70 change like “has a purple star pattern”, while the latter emphasizes more

71 overall change such as “is modern and fashionable”. • Unlike the image-text
72 compositor that uses relative caption to describe **what to change** on the
73 reference image, we devise the text-image compositor to capture the tex-
74 tual cues based on text-to-image retrieval, where the reference image implies
75 **what to preserve** for the reference image. See details in Sec. 3.1. To min-
76 imize the negative impact of triplet ambiguity during training, we impose
77 a KL loss between two image-text compositors. The KL loss promotes two
78 compositors to learn from each other and reach a consensus, which is similar
79 to supervision from peers in a group, as it helps each compositor to reduce
80 its own bias and thus avoids overfitting to the annotated target image.

81 In summary, our contributions are as follows:

82 • We have identified an inherent issue within the context of composed
83 image retrieval, namely triplet ambiguity, which we subsequently confirm
84 through initial experimental investigations (*see Fig. 2 and 4*). This problem,
85 stemming from the inherent noisiness of the annotation process, results in
86 suboptimal model learning, as it compromises the extraction of discriminative
87 features that integrate visual and linguistic information.

88 • To relieve triplet ambiguity, we introduce the Consensus Network (Css-
89 Net) featuring a consensus module with four distinct compositors for collab-
90 orative training (*see Table 5*) and joint inference (*see Table 6*).

91 • Extensive experiments show that the proposed method minimizes the
92 negative impacts of noisy triplets. On three prevalent public benchmarks,
93 we observe that Css-Net significantly surpasses the current state-of-the-art
94 competitive methods, *e.g.*, with +2.77% Recall@10 on Shoes, and +6.67%
95 Recall@50 on FashionIQ (*see Table 1, 2, and 3*).

96 2. Related Work

97 **Cross-modal Image Retrieval.** Cross-modal image retrieval has attracted
98 wide attention from researchers. The most popular patterns of image re-
99 trieval are image-to-image matching (Zheng et al., 2017; Deng et al., 2019;
100 Sun et al., 2020; Wu et al., 2017; Dai et al., 2018; Liu et al., 2022; Qu et al.,
101 2024) and text-to-image matching (Liu et al., 2019; Zhang et al., 2020; Liu
102 et al., 2022; Zhang et al., 2022; Li et al., 2024). Although these paradigms
103 have made great progress, they do not provide enough convenience for users
104 to express their search intention. Therefore, more forms of image retrieval
105 with flexible queries such as sketch-based image retrieval (Deng et al., 2020;
106 Wang et al., 2021; Li et al., 2022; Liang et al., 2024) have emerged. In this

107 work, the composed image retrieval task involves a composed query of a ref-
108 erence image and a relative caption. To tackle this task, recent works (Vo
109 et al., 2019; Chen et al., 2020; Yang et al., 2021; Zhang et al., 2021; Lee et al.,
110 2021; Wen et al., 2021; Gu et al., 2021; Zhao et al., 2022; Han et al., 2023)
111 devise diverse composition architectures to capture the visual-linguistic re-
112 lation. Unlike the methods described above, our Css-Net does not propose
113 complicated compositors. Instead, our work mainly focuses on reducing sin-
114 gular compositor biases to alleviate the identified triplet ambiguity problem.

115 **Attention Mechanism.** The attention mechanism is widely used in lan-
116 guage and vision tasks in machine learning to capture the relations between
117 features. This mechanism is also inspired by a psychological finding (Cor-
118 betta and Shulman, 2002) that humans observe and pay attention to specific
119 parts as needed. In the composed image retrieval task, many works use
120 the attention mechanism to design the image-text compositor. For example,
121 VAL (Chen et al., 2020) employs self-attention to capture the image-text re-
122 lations by concatenating the text feature to the image feature. CoSMo (Lee
123 et al., 2021) adopts the disentangled multi-modal non-local block to stabilize
124 the training procedure for learning better representations. Besides, CLVC-
125 Net (Wen et al., 2021) proposes a cross attention between each word in the
126 sentence and each spatial location of the image feature to recognize details.
127 In our work, the main idea is not to design a new attention-based compositor
128 but to utilize several compositors to form as a consensus module. Without
129 loss of generality, we deploy the widely-used CoSMo as the image-text com-
130 positor. Moreover, we propose specific text-image compositors based on cross
131 attention to better capture the relation between the reference image feature
132 and the word-level text feature, which is orthogonal with existing attention-
133 based models and could further improve the retrieval performance.

134 **Co-training.** Co-training is a semi-supervised learning technique that ex-
135 ploits two components to acquire complementary information on two views
136 of the data (Blum and Mitchell, 1998). It has been extensively utilized in
137 various research fields such as image recognition (Qiao et al., 2018), segmen-
138 tation (Peng et al., 2020; Hui et al., 2023) and domain adaptation (Saito et al.,
139 2018; Zheng and Yang, 2019; Luo et al., 2019). Our work adopts a co-training
140 paradigm that leverages four compositors with different knowledge to jointly
141 make decisions for the composed image retrieval task. The two image-text
142 compositors focus on the detailed and overall changes to the reference images
143 based on the perspective of finding “what to change” in the reference image,
144 and the two text-image compositors are in view of the text-to-image retrieval

145 with the reference image implying “what to preserve” for the relative cap-
 146 tion. The compositors hold diverse knowledge from different views of the
 147 data. Thus, we explicitly encourage the consensus between compositors and
 148 leverage the consensus to rectify the single prediction.

149 3. Method

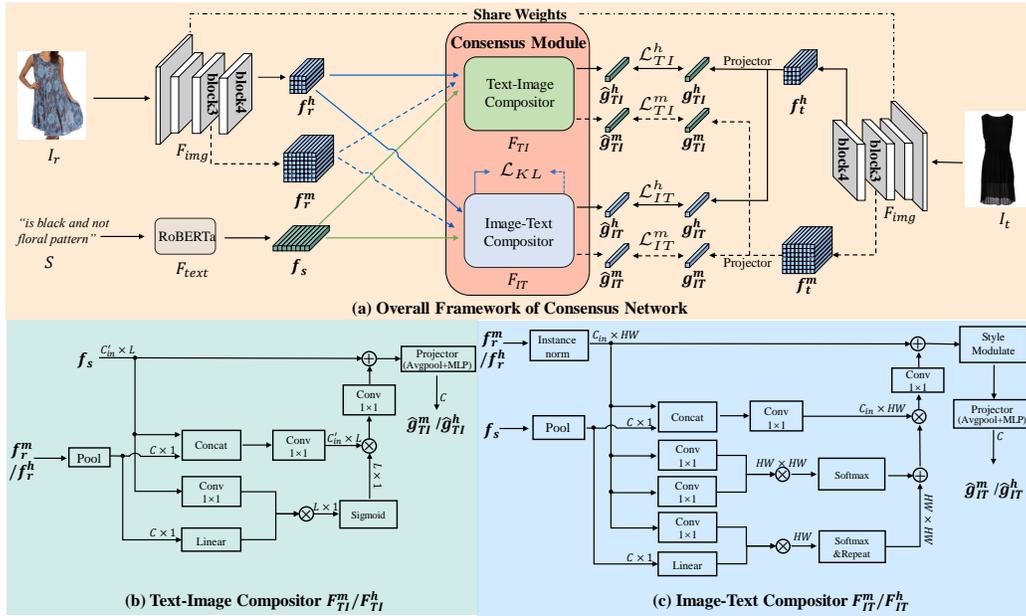


Figure 3: Schematic illustration of the Consensus Network. Given a reference image and a relative caption, the image encoder F_{img} extracts the mid-level image feature f_r^m and high-level image feature f_r^h , and the text encoder F_{text} extracts the text feature f_s . Then, compositors fuse the text feature with either the mid-level or high-level image feature. Each compositor generates distinct composed feature. Finally, we match the composed features with the corresponding target features and impose a KL loss between image-text compositors for training.

150 3.1. Overview of Consensus Network

151 As illustrated in Fig. 3 (a), the Consensus Network consists of three
 152 components: a image encoder, a text encoder, and a consensus module.
 153 The image encoder, F_{img} , extracts mid-level and high-level reference im-
 154 age features as $f_r^m, f_r^h = F_{img}(I_r)$, where I_r is the reference image, and

155 $\mathbf{f}_r^m, \mathbf{f}_r^h \in \mathbb{R}^{C_{in} \times (H \times W)}$ are mid-level and high-level image features, respec-
 156 tively (*i.e.*, output from block3 and block4 of the ResNet (He et al., 2016)).
 157 $C_{in} \times (H \times W)$ represents the shape of the feature maps. For brevity, we do
 158 not distinguish between different shapes of image features. The text encoder,
 159 denoted as F_{text} , extracts features of the relative caption as $\mathbf{f}_s = F_{text}(S)$,
 160 where S denotes the relative caption, $\mathbf{f}_s \in \mathbb{R}^{C'_{in} \times L}$ refers to the word-level
 161 representation, and L is the number of words of the relative caption.

162 After extracting the image and text features, the consensus module trans-
 163 forms the reference image features with the corresponding text features into
 164 the composed features. It consists of four distinct compositors possessing
 165 different knowledge. These compositors at different depths of the image en-
 166 coder can be grouped into two types. Specifically, given the reference image
 167 feature \mathbf{f}_r and the text feature \mathbf{f}_s , the composed query $\hat{\mathbf{g}}$ can be obtained by
 168 either an image-text compositor or a text-image compositor. The image-text
 169 compositor has the residual form of $\hat{\mathbf{g}}_{IT} = \mathbf{f}_r + comp(\mathbf{f}_r, \mathbf{f}_s)$ and mainly fo-
 170 cuses on “what to change” for \mathbf{f}_r conditioned on the relative caption, while
 171 the text-image compositor has the residual form of $\hat{\mathbf{g}}_{TI} = \mathbf{f}_s + comp(\mathbf{f}_s, \mathbf{f}_r)$
 172 and mainly emphasizes on “what to preserve” for \mathbf{f}_s conditioned on the re-
 173 ference image. Here, $comp$ represents a trained component to fuse \mathbf{f}_r and \mathbf{f}_s
 174 as the condition. Considering both the performance and computational effi-
 175 ciency, the text-image compositors F_{TI}^m and F_{TI}^h , shown in Fig. 3 (b), take the
 176 word-level representation \mathbf{f}_s along with the average pooled reference image
 177 features $pool(\mathbf{f}_r^m), pool(\mathbf{f}_r^h)$ as input, respectively:

$$\left\{ \begin{array}{l} \hat{\mathbf{g}}_{TI}^m = F_{TI}^m(\mathbf{f}_s, pool(\mathbf{f}_r^m)) \\ \hat{\mathbf{g}}_{TI}^h = F_{TI}^h(\mathbf{f}_s, pool(\mathbf{f}_r^h)), \end{array} \right. \quad (1)$$

178 where $\hat{\mathbf{g}}_{TI}^m, \hat{\mathbf{g}}_{TI}^h$ are the composed features from text-image compositors.
 179 Similarly, the image-text compositors F_{IT}^m, F_{IT}^h , shown in Fig. 3 (c) take the
 180 intermediate image feature maps, $\mathbf{f}_r^m, \mathbf{f}_r^h$ along with the pooled sentence-
 181 level text representation $pool(\mathbf{f}_s)$ as input, which are given by:

$$\left\{ \begin{array}{l} \hat{\mathbf{g}}_{IT}^m = F_{IT}^m(\mathbf{f}_r^m, pool(\mathbf{f}_s)) \\ \hat{\mathbf{g}}_{IT}^h = F_{IT}^h(\mathbf{f}_r^h, pool(\mathbf{f}_s)), \end{array} \right. \quad (2)$$

182 where $\hat{\mathbf{g}}_{IT}^m, \hat{\mathbf{g}}_{IT}^h$ are the composed features from image-text compositors.

183 The target image features $\mathbf{f}_t^m, \mathbf{f}_t^h$ are obtained from the same image
 184 encoder F_{img} as the reference image features $\mathbf{f}_r^m, \mathbf{f}_r^h$. Then four independent

185 projector blocks (composed of an average pooling layer and a MLP) are
 186 employed to acquire target features: \mathbf{g}_{TI}^m , \mathbf{g}_{TI}^h , \mathbf{g}_{IT}^m , and \mathbf{g}_{IT}^h . Finally, the
 187 four compositors are trained by pulling close the corresponding target while
 188 pushing away other negatives within the embedding space.

189 3.2. Consensus Module

190 To relieve the triplet ambiguity, we introduce the consensus module,
 191 which consists of four distinct compositors with different knowledge. These
 192 compositors have individual biases learned on noisy triplets, which are mini-
 193 mized at two stages. At the training stage, each compositor acquires informa-
 194 tion from different views of the data, and the KL loss enables them to learn
 195 from each other to minimize biases. At the evaluation stage, each compositor
 196 independently provides decisions and collaborates to rank the entire gallery
 197 by aggregating their decisions. We first discuss the design to ensure that
 198 each compositor acquires distinct knowledge, then explain how compositors
 199 learn from each other to reduce their biases learned on noisy triplets. The
 200 batch-based classification loss is as follows:

$$\mathcal{L}_{BBC} = -\log \frac{\exp(\hat{\mathbf{g}} \cdot \mathbf{g}_+)}{\sum_{j=1}^B \exp(\hat{\mathbf{g}} \cdot \mathbf{g}_j)}, \quad (3)$$

201 where $\hat{\mathbf{g}}$ is the composed feature from respective compositor, and \mathbf{g}_j are
 202 candidates, among which the true match is \mathbf{g}_+ .

203 **Pyramid Training for Image-Text Compositor.** We develop a pyra-
 204 mid training paradigm for image-text compositors, which is inspired by the
 205 finding (Lin et al., 2017; Miech et al., 2021) that the image features of high-
 206 resolution are semantically weak, while the image features of low-resolution
 207 are semantically strong. Through exploring the different spatial information
 208 of the reference image, the two image-text compositors F_{IT}^m and F_{IT}^h inde-
 209 pendently learn knowledge by leveraging the batch-based classification loss
 210 \mathcal{L}_{IT}^m and \mathcal{L}_{IT}^h . The independent batch-based classification loss makes each
 211 image-text compositor learn from the interactions between relative caption
 212 and different spatial information of the reference image, which enables these
 213 compositors to hold distinct knowledge from each other.

214 **Auxiliary knowledge from Text-Image Compositor.** The text-image
 215 compositor is a fancy component for generating the composed feature from
 216 the input, which is seldom referred to in previous works. It offers additional
 217 knowledge due to its distinct design from the image-text compositor. As

218 discussed in Sec. 3.1, the text-image compositor views the data from another
 219 perspective, mainly focusing on the text-to-image retrieval with the reference
 220 image implying “what to preserve” conditioned on the text information, while
 221 the image-text compositor finds “what to change” in the reference image. We
 222 use two symmetric text-image compositors at the same depths of the image
 223 encoder, leveraging the batch-based classification loss \mathcal{L}_{TI}^m and \mathcal{L}_{TI}^h .

224 **Collaborative Consensus Learning.** The triplet ambiguity problem leads
 225 to noisy triplets and biases the model learning. To mitigate this problem,
 226 we use the Kullback Leibler divergence loss (KL loss) for two image-text
 227 compositors. The KL loss promotes the compositors to learn from each
 228 other, reducing biases and reaching a consensus. This approach balances
 229 the preservation of distinct knowledge and the attainment of consensus. By
 230 enhancing cooperation and knowledge sharing, our method is more robust to
 231 the triplet ambiguity problem. Specifically, we denote the resulting posterior
 232 probability of F_{IT}^m as \mathbf{p}^m and that of F_{IT}^h as \mathbf{p}^h . We set a target probability
 233 \mathbf{p}^w as the weighted sum of both \mathbf{p}^m and \mathbf{p}^h , which is given by:

$$\mathbf{p}^w = \lambda_1 \cdot \mathbf{p}^m + \lambda_2 \cdot \mathbf{p}^h, \quad (4)$$

234 where λ_1 and λ_2 are weight coefficients, and the KL loss is formulated as:

$$\mathcal{L}_{KL} = D_{KL}(\mathbf{p}^m || \mathbf{p}^w) + D_{KL}(\mathbf{p}^h || \mathbf{p}^w), \quad (5)$$

235 where D_{KL} is the KL divergence distance. The KL loss reduces the biases
 236 of the compositors during training, which works alongside the batch-based
 237 classification loss in our approach. The preliminary experiments show that
 238 it is not essential to incorporate extra KL loss for the two text-image com-
 239 positors. See Sec. 5.3 for a detailed explanation. The final loss for training
 240 is the sum of the above loss functions:

$$\mathcal{L} = \mathcal{L}_{IT}^m + \mathcal{L}_{IT}^h + \mathcal{L}_{TI}^m + \mathcal{L}_{TI}^h + \mathcal{L}_{KL}, \quad (6)$$

241 where L_{IT}^m , L_{IT}^h , L_{TI}^m , and L_{TI}^h are batch-based classification loss used for in-
 242 dependently training each image-text/text-image compositor F_{IT}^m , F_{IT}^h , F_{TI}^m ,
 243 and F_{TI}^h . The superscript m indicates the mid-level input feature, while h
 244 denotes the high-level feature. The subscript IT indicates the image-text
 245 compositor, while the subscript TI denotes the text-image compositor.

246 **Joint Inference.** The four distinct compositors independently learn differ-
 247 ent knowledge from the data triplets and enable the knowledge transfer to re-
 248 duce biases learned on noisy triplets. At the evaluation step, we involve each

249 compositor in decision-making to further minimize individual bias. Specifi-
 250 cally, we use each compositor to independently generate composed features
 251 and measure the similarity between any composed feature and target feature.
 252 The resulting similarity matrices are denoted as $P_{IT}^m, P_{IT}^h, P_{TI}^m, P_{TI}^h \in \mathbb{R}^{n_1 \times n_2}$,
 253 where n_1 and n_2 are the number of queries and target images in the gallery.
 254 The final similarity matrix for ranking the gallery is the weighted sum of four
 255 similarity matrices from distinct compositors:

$$P = \alpha_1 \cdot P_{IT}^m + \alpha_2 \cdot P_{IT}^h + \alpha_3 \cdot P_{TI}^m + \alpha_4 \cdot P_{TI}^h, \quad (7)$$

256 where $\alpha_1 \dots \alpha_4$ are weight coefficients to balance the decisions from four com-
 257 positors. Note that a common practice that concatenates multiple composed
 258 features as one query is a special case where all α s are equal to 1.

259 4. Experiments

260 4.1. Experimental Setup

261 **Datasets.** We evaluate Css-Net on three composed image retrieval datasets,
 262 *i.e.*, Shoes (Berg et al., 2010), FashionIQ (Wu et al., 2021), and Fashion200k (Vo
 263 et al., 2019).

- 264 • The Shoes dataset (Berg et al., 2010) is originally crawled from *like.com*
 265 for attribute discovery. It is then annotated in the form of a triplet for
 266 dialog-based interactive retrieval. We follow VAL (Chen et al., 2020)
 267 to use 10,000 training samples and 4,658 evaluation samples.
- 268 • The FashionIQ dataset (Wu et al., 2021) is a language-based interactive
 269 fashion retrieval dataset with 77,684 images across three categories:
 270 Dresses, Tops&Tees, and Shirts. It includes 18,000 triplets from 46,609
 271 training images, each containing a reference image, a target image, and
 272 two descriptive natural language captions. The evaluation procedure
 273 follows VAL (Chen et al., 2020) and CoSMo (Lee et al., 2021).
- 274 • The Fashion200k dataset (Han et al., 2017) contains over 200k fash-
 275 ion images from various websites and is for attribute-based product
 276 retrieval. With descriptive attributes for each product, 172k images
 277 are used for training and 33,480 test queries for evaluation, following
 278 VAL and CoSMo methods. The relative descriptions are generated
 279 from attributes using an online-processing pattern.

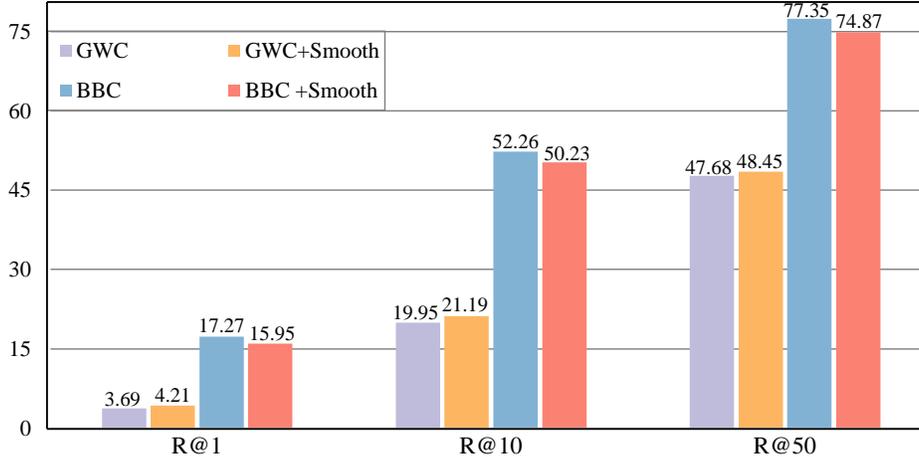


Figure 4: Comparison between the batch-based classification and the global-wise classification (GWC) on the Shoes dataset. GWC significantly degrades the performance since more false negative samples are involved due to triplet ambiguity.

280 *4.2. Triplet Ambiguity Verification*

281 **Global-wise *v.s.* Batch-based Optimization.** To verify the negative
 282 impacts from the noisy triplets as shown in Fig. 2, we quantitatively com-
 283 pare global-wise with batch-based optimization objectives. In particular, •
 284 Batch-based Classification (BBC): Limited negatives in the current batch
 285 are involved, and • Global-wise Classification (GWC): Mining more negative
 286 samples in the whole training set for comparison.

287 If the data triplets do **NOT** have ambiguity, the global-wise classification
 288 has the potential to be comparable or even better since it uses more negative
 289 samples in the training set and potentially learns a better metric, which is
 290 consistent with many findings in metric learning (Hermans et al., 2017; Sheng
 291 et al., 2020; Wang et al., 2020) and self-supervised learning (Chen et al.,
 292 2020; He et al., 2020). Specifically, Given a query q and features/prototypes
 293 $\{k_0, k_1, \dots\}$ of candidate target images, where the true match is denoted as
 294 k_+ . Two losses are given by:

$$\mathcal{L}_{BBC} = -\log \frac{\exp(q \cdot k_+)}{\sum_{i=1}^B \exp(q \cdot k_i)} \quad (8)$$

295 and

$$\mathcal{L}_{GWC} = -\log \frac{\exp(q \cdot k_+)}{\sum_{i=1}^N \exp(q \cdot k_i)}, \quad (9)$$

Method	Dress		Shirt		Toptee		Average	
	R@10 ↑	R@50 ↑						
MRN (Kim et al., 2016)	12.32	32.18	15.88	34.33	18.11	36.33	15.44	34.28
FiLM (Perez et al., 2018)	14.23	33.34	15.04	34.09	17.30	37.68	15.52	35.04
TIRG (Vo et al., 2019)	14.87	34.66	18.26	37.89	19.08	39.62	17.40	37.39
VAL (Chen et al., 2020)	21.12	42.19	21.03	43.44	25.64	49.49	22.60	45.04
DCNet (Kim et al., 2021)	28.95	56.07	23.95	47.30	30.44	58.29	27.78	53.89
CoSMo* (Lee et al., 2021)	26.45	52.43	26.94	52.99	31.95	62.09	28.45	55.84
CLVC-Net† (Wen et al., 2021)	29.85	56.47	28.75	54.76	33.50	64.00	30.70	58.41
ARTEMIS (Delmas et al., 2022)	27.16	52.40	21.78	54.83	29.20	43.64	26.05	50.29
MUR (Chen et al., 2022)	30.60	57.46	31.54	58.29	37.37	68.41	33.17	61.39
CLIP4Cir (Baldrati et al., 2022)	31.73	56.02	35.77	57.02	36.46	62.77	34.65	58.60
Baseline	30.95	56.98	31.48	59.98	36.97	67.31	33.13	61.42
Css-Net	33.65	63.16	35.96	61.96	42.65	70.70	37.42	65.27

Table 1: Quantitative results on the FashionIQ dataset. The best results are in **bold**. The symbol * marks an updated results by the same authors. The symbol † indicates that this method deploys model ensemble (the same as below).

296 where B is the batch size, and N is the number of IDs (classes) in the training
 297 set. The only difference between them is that \mathcal{L}_{GWC} involves more negative
 298 counterparts, which results in high false negative rates if the triplet ambigu-
 299 ity does exist. We conduct experiments on the Shoes dataset (Berg et al.,
 300 2010) using two losses, respectively, under the same settings of CoSMo (Lee
 301 et al., 2021). We observe that batch-based methods outperform global-wise
 302 methods by a large margin, as shown in Fig. 4. The experimental results
 303 confirm our triplet ambiguity assumption: the training data contains many
 304 noisy triplets (*i.e.*, false negative samples). Although batch-based classifica-
 305 tion suffers less from triplet ambiguity, the single compositor still faces some
 306 noisy negative triplets in the batch and produces a sub-optimal solution.

307 **Label Smoothing.** One intuitive way we consider to alleviate the triplet
 308 ambiguity problem is label smoothing. The motivation is that there are
 309 many false negative samples due to the triplet ambiguity problem, and label
 310 smoothing could alleviate the overfitting to the annotated true match. In
 311 label smoothing, the label $\mathbf{y} = [y_1, \dots, y_n]$ is not a hard one-hot label rather
 312 than a soft one-hot label, which is given by:

$$y_i = \begin{cases} 1 & (\text{if } i = c) \\ 0 & (\text{if } i \neq c) \end{cases} \implies y_i = \begin{cases} 1 - \varepsilon & (\text{if } i = c) \\ \frac{\varepsilon}{B-1} & (\text{if } i \neq c), \end{cases} \quad (10)$$

313 where y_i is the label for class i , c is the corresponding class of the query, B

Method	Shoes		
	R@1 \uparrow	R@10 \uparrow	R@50 \uparrow
MRN (Kim et al., 2016)	11.74	41.70	67.01
FiLM (Perez et al., 2018)	10.19	38.89	68.30
TIRG (Vo et al., 2019)	12.60	45.45	69.39
VAL (Chen et al., 2020)	16.49	49.12	73.53
CoSMo (Lee et al., 2021)	16.72	48.36	75.64
DCNet (Kim et al., 2021)	-	53.82	79.33
CLVC-Net \dagger (Wen et al., 2021)	17.64	54.39	79.47
MUR (Chen et al., 2022)	18.41	53.63	79.84
ARTEMIS (Delmas et al., 2022)	18.72	53.11	79.31
Baseline	17.27	52.26	77.35
Css-Net	20.13	56.81	81.32

Table 2: Quantitative results on the Shoes dataset. The best results are in **bold**. The symbol \dagger indicates that this method deploys model ensemble. The proposed method has achieved competitive performances in all three metrics R@1, 10, 50.

314 is the batch size, and ε is a hyperparameter for label smoothing and is set
315 to be 0.1. We use label smoothing for both the batch-based classification
316 and the global-wise classification, which are presented in Fig. 4. The experi-
317 mental results indicate that label smoothing deteriorates the performance of
318 batch-based classification but enhances the performance of global-wise clas-
319 sification. This is because • global-wise classification is severely affected by
320 triplet ambiguity due to high false negative rate, while batch-based classifi-
321 cation is affected only when noisy negative triplets are in the batch; • Label
322 smoothing could relieve the triplet ambiguity but introduce another problem
323 that many true negative target samples are assigned weights, which impairs
324 the model training for batch-based classification. The experimental results
325 also verify the effectiveness of KL loss as another form of soft label.

326 4.3. The Effectiveness of Our Method

327 We present the experimental results in Table 1, Table 2, and Table 3.
328 We could make two observations: **(1) We adopt a competitive base-**
329 **line with few modifications.** As mentioned in Sec. 4.1, we adopt the
330 CoSMo as our baseline and replace the LSTM with a more robust text en-
331 coder: RoBERTa, and observe consistent improvement. For example, on the
332 FashionIQ dataset, our baseline improves CoSMo by 4.68% R@10 on average
333 and surpasses CoSMo by 3.90% R@10 on the Shoes dataset. We infer that

Method	Fashion200k		
	R@1 \uparrow	R@10 \uparrow	R@50 \uparrow
MRN (Kim et al., 2016)	13.4	40.0	61.9
FiLM (Perez et al., 2018)	12.9	39.5	61.9
TIRG (Vo et al., 2019)	14.1	42.5	63.8
VAL (Chen et al., 2020)	21.2	49	68.8
DCNet (Kim et al., 2021)	-	46.9	67.6
CoSMo (Lee et al., 2021)	23.3	50.4	69.3
CLVC-Net \dagger (Wen et al., 2021)	22.6	53.0	72.2
ARTEMIS (Delmas et al., 2022)	21.5	51.1	70.5
Baseline	20.9	47.7	67.8
Css-Net	22.2	50.5	69.7
Css-Net \dagger	23.4	52.0	72.0

Table 3: Quantitative results on the Fashion200k dataset. The best results are in **bold**. The symbol \dagger indicates that this method deploys model ensemble. The proposed method has achieved competitive performances.

334 RoBERTa is more robust than LSTM (Hochreiter and Schmidhuber, 1997) to
335 accurately capture the textual information. However, our baseline is slightly
336 lower than the reported results of CoSMo on Fashion200k, as the authors do
337 not provide sufficient implementation details for reproducing. This also lim-
338 its comparing our method with CQBIR (Zhang et al., 2022), whose baseline
339 uses faster RCNN (Girshick, 2015) as a different image encoder. Nevertheless,
340 our method is more effective than CQBIR on FashionIQ and Shoes, where
341 the triplet ambiguity problem is more serious. **(2) The proposed Css-Net**
342 **could further improve and advances the state of the art on such a**
343 **strong baseline, verifying the effectiveness of Css-Net.** For example,
344 Table 1 shows Css-Net improves retrieval accuracy on all FashionIQ subsets.
345 Compared to the baseline, it gains +2.70% R@10 on Dress, +4.48% R@10
346 on Shirt, and +5.68% R@10 on TopTee. Compared to previous works, our
347 method brings overall improvements (e.g., +2.77% R@10 and +6.67% R@50
348 on average by CLIP4Cir). The improvements are significant and empirically
349 validate the effectiveness of Css-Net for handling the triplet ambiguity prob-
350 lem. Besides in Table 2, Css-Net surpasses the state-of-the-art (CLVC-Net)
351 on the Shoes dataset, achieving improvements of +2.49% R@1 and +2.42%
352 R@10, which further demonstrates that Css-Net is robust across different
353 datasets. Table 3 presents Fashion200k results. Although our baseline is

Method	Shoes		
	R@1 ↑	R@10 ↑	R@50 ↑
F_{IT}^h	17.27	52.26	77.35
$F_{IT}^l + F_{IT}^h$	18.24	52.14	78.12
$F_{IT}^l + F_{IT}^m + F_{IT}^h$	18.81	54.21	79.55
$F_{IT}^m + F_{IT}^h$	19.10	54.69	79.63

Table 4: Comparison of various pyramid training methods on the Shoes dataset. These methods are trained and evaluated independently. F_{IT}^l , F_{IT}^m , and F_{IT}^h represent the low-level, mid-level, and high-level image-text compositor, respectively. The low-level compositor is useful, whereas the mid and high-level features show better performance.

\mathcal{L}_{IT}^m	$\mathcal{L}_{TI}^h + L_{TI}^m$	\mathcal{L}_{KL}	Shoes		
			R@1 ↑	R@10 ↑	R@50 ↑
Baseline: (only \mathcal{L}_{IT}^h)			17.27	52.26	77.35
✓			19.10(+1.83)	54.69(+2.43)	79.63(+2.28)
✓	✓		19.47(+2.20)	54.63(+2.37)	80.46(+3.11)
✓	✓	✓	20.13(+2.86)	56.81(+4.55)	81.32(+3.97)

Table 5: Efficacy of model designs. L_{IT}^m , L_{IT}^h , L_{TI}^m , and L_{TI}^h are batch-based classification loss defined in Eqn. 3, and L_{KL} is the KL loss defined in Eqn. 5.

354 below the reported results of CosMo because of insufficient implementation
355 details for reproduction, Css-Net brings a considerable improvement (*e.g.*,
356 +2.8% R@10 over the baseline) and is still competitive with many SOTA
357 works especially when applying the model ensemble (*e.g.*, +4.3% R@10).

358 4.4. Diagnostic Experiments

359 **Pyramid Training.** In Sec. 3.2, we present the design of the pyramid train-
360 ing, which exploits the image features from the mid-level and high-level blocks
361 of the image encoder. We verify its effectiveness by comparing it with dif-
362 ferent designs. Table 4 reports the experimental results. Our baseline is
363 $F_{IT}^m + F_{IT}^h$ used in Css-Net. We conduct experiments on two variants for
364 pyramid training: 1) $F_{IT}^l + F_{IT}^h$, which uses the image features from block2
365 and block4 of the ResNet, and 2) $F_{IT}^l + F_{IT}^m + F_{IT}^h$, utilizing three image-text
366 compositors at three depths. Both variants perform worse than Css-Net,
367 *e.g.*, -2.55% and -0.48% on the R@10 metric. However, they both surpass

Inference Method	Shoes		
	R@1 ↑	R@10 ↑	R@50 ↑
F_{IT}^m	15.72	51.17	78.89
F_{IT}^h	18.35	55.15	80.52
F_{TI}^m	17.06	53.35	78.92
F_{TI}^h	16.58	52.17	77.77
Joint Inference (Eq. 7)	20.13	56.81	81.32

Table 6: Effect of joint inference. We train Css-Net with four compositors on Shoes once and separately evaluate each compositor. Joint inference refers to using the weighting scheme (Eqn. 7) to combine decisions from all the compositors .

	Total time (s) ↓	Time per query (ms) ↓	Time per target (ms) ↓
Baseline (one)	168.2	50.2	56.4
Css-Net (four)	195.8	58.5	65.7

Table 7: Inference time cost for the baseline and Css-Net. Total time refers to the time taken to process all queries. Time per query indicates the average time spent on each query, while time per target represents the average time used to process each target in the gallery.

368 F_{IT}^h using only one image-text compositor at block4. These results indicate
369 that 1) the low-level image feature is too semantically weak to provide image
370 information, and 2) groups perform better than individuals.

371 **Efficacy of Model Designs.** Table 5 shows the effectiveness of our core
372 idea, which uses four different compositors with KL loss to relieve the triplet
373 ambiguity problem. We make three observations from the table. First, em-
374 ploying image-text compositors at other layers of the image encoder (*i.e.*,
375 \mathcal{L}_{IT}^m) can mitigate the triplet ambiguity problem and improve the perfor-
376 mance significantly (77.35% \rightarrow 79.63% at R@50 metric). This indicates that
377 two image-text compositors can benefit from the interactions between the
378 relative caption and different spatial information of the reference image. Sec-
379 ond, adding a new compositor module, text-image compositor, to this task
380 (*i.e.*, $\mathcal{L}_{TI}^m + \mathcal{L}_{TI}^h$) can further improve the performance (79.63% \rightarrow 80.46%
381 at R@50 metric). This demonstrates the advantage of auxiliary knowledge.
382 Third, applying an extra KL loss for two image-text compositors (\mathcal{L}_{KL}) can
383 enhance the performance notably (80.46% \rightarrow 81.32% at R@50 metric). This
384 suggests that the KL loss enables two image-text compositors to share their

385 knowledge, thus minimizing the biases learned from noisy triplets.

386 **Effect of Joint Inference** At the evaluation stage, Css-Net allows com-
387 positors to jointly make the decision as introduced in Sec. 3.2. As shown in
388 Table 6, joint inference surpasses single compositor and verifies our motiva-
389 tion that groups perform better than individuals and could be used to reduce
390 their own prediction biases mainly caused by the triplet ambiguity problem.

391 **Computational cost at inference** Css-Net uses four compositors that
392 share the same image and text encoders, thus adding minimal retrieval la-
393 tency. The inference time is shown in Table 7, ranging from loading the
394 model to displaying results. The experiments are conducted with GeForce
395 RTX 2080 Ti, using 33,480 queries and 29,789 targets.

396 **Implementation Details.** We modify CoSMo (Lee et al., 2021) as our
397 baseline by replacing LSTM (Graves, 2012) with RoBERTa (Liu et al., 2019)
398 as the text encoder. ResNet-50 (He et al., 2016) serves as the image encoder
399 for Shoes and FashionIQ datasets, while ResNet-18 (He et al., 2016) is used
400 for Fashion200k. Embedding space dimension C is 512. Text feature shape is
401 $C'_{in} \times L$, with C'_{in} being 768 and L is the sentence length. During training, we
402 set $\lambda_1 = 10$ and $\lambda_2 = 1$, while evaluation uses $\alpha_1 \dots \alpha_4 = 1, 0.5, 0.5, 0.5$. We
403 adopt the standard evaluation metric in retrieval, *i.e.*, Recall@K, denoted as
404 R@K for short. We use a random seed for each experiment and repeat it five
405 times for the final results. we employ the Adam optimizer (Kingma and Ba,
406 2014) with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. On Shoes and FashionIQ, the batch
407 size is set to be 32 and the base learning rates of the text encoder and other
408 modules are $2e - 6$ and $2e - 5$, respectively. On Fashion200k, the batch size
409 is set to be 128 and the base learning rates are $2e - 6$ and $2e - 4$, respectively.
410 We adopt warm-up for the first 5 epochs, decay learning rate by 10 at epochs
411 35 and 45 during training. The total training epoch is 50.

412 5. Further Analysis and Discussion

413 5.1. Further Qualitative Analysis

414 Fig. 5 shows the top-10 retrieval results on three datasets: Shoes, Fash-
415 ion200K, and FashionIQ. We make three key observations from these results:
416 (1) Css-Net can capture the information of the reference image and the rel-
417 ative caption for both coarse-grained and fine-grained queries. For example,
418 the first query of Shoes and the third query of FashionIQ retrieve the cor-
419 rect matches easily, and the first and second queries of FashionIQ also find
420 the correct matches. These queries have clear and distinctive features that

	Composed Query	Top 10 Retrieval Results Rank 1→10
Shoes	 are red	         
	 is black not brown	         
	 have black and grey plaid on the middle	         
FashionIQ	 is orange with black lettering	         
	 longer sleeves and is brown	         
	 is yellow	         
Fashion200k	 replace black with white	         
	 replace yellow with blue	         
	 replace black with white	         

Figure 5: Top-10 retrieval results on three datasets. The composed queries consist of a reference image and a relative caption that describes the desired modification. The blue/green boxes refer to the reference image and the true match(es).

421 can be matched by Css-Net. (2) The model sometimes fails to retrieve the
422 correct matches due to the triplet ambiguity problem, *e.g.*, the first query
423 of Fashion200K retrieves some negative samples but are still highly related
424 to the query. (3) Css-Net is less sensitive to some detailed information such
425 as location. For example, the third query in Shoes retrieves a shoe that
426 is visually similar but has a wrong paid location, because the dataset has
427 few similar training samples. Improving the sensitivity of the model to the
428 detailed information is a direction for our future work. We plan to explore
429 more fine-grained features to enhance Css-Net in the future works.

430 *5.2. Comparison with Most Relevant Works*

431 We compare our *Css-Net* with VAL (Chen et al., 2020) and CLVC-Net (Wen
432 et al., 2021), which are most relevant to our work.

433 (1) Our *Css-Net* differs from the hierarchical matching strategies in VAL:

434 • Our *Css-Net* facilitates knowledge sharing between composers at var-
435 ious depths for consensus, instead of independent learning in VAL.

436 • Our *Css-Net* observes that the low-level composer does not contribute
437 to collaborative learning and omits it to enhance the recall performance
438 and efficiency (*See Table 1*).

439 • Our *Css-Net* implements an adjustable weighted sum during evaluation,
440 enabling individuals to make decisions as a group.

441 (2) Our *Css-Net* differs from the model ensemble design in CLVC-Net:

442 • Our *Css-Net* is more efficient since all composers share the same en-
443 coder stem (Table 7), while model ensembling in CLVCNet employs
444 several independent backbones.

445 • Our *Css-Net* encourages intra-modal and inter-modal knowledge shar-
446 ing via collaborative learning between composers, while model ensem-
447 bling does not entail additional loss or learning among the models.

448 • Our *Css-Net* acknowledges that the composers have different knowl-
449 edge and thus assign adaptive weights, while model ensembling usually
450 presumes that the models are independent and equally important.

451 • Our *Css-Net* enables single composer to perform better could further
452 benefits from model ensembling (*See Table 6*), while model ensembling
453 does not improve single composer prediction.

454 *5.3. Discussion of Collaborative Learning*

455 We apply a KL loss between text-image composers in a preliminary ex-
456 periment, but find that it is not as significant as the KL loss between image-
457 text composers. This is because the inputs for the text-image composers
458 are too similar, as shown in Fig. 6. Specifically, both text-image composi-
459 tors receive a pooled reference image feature with identical dimensions and
460 share the same text representations. Therefore, the main function of these

461 text-image compositors is to act as auxiliary decision-makers during joint in-
 462 ference, addressing the triplet ambiguity issue. For simplicity and efficiency,
 463 we do not incorporate additional KL loss for the text-image compositors.
 464 However, we note that the text-image compositors still play an important
 465 role in our framework, as they provide complementary information to the
 466 image-text compositors and improve the retrieval performance.

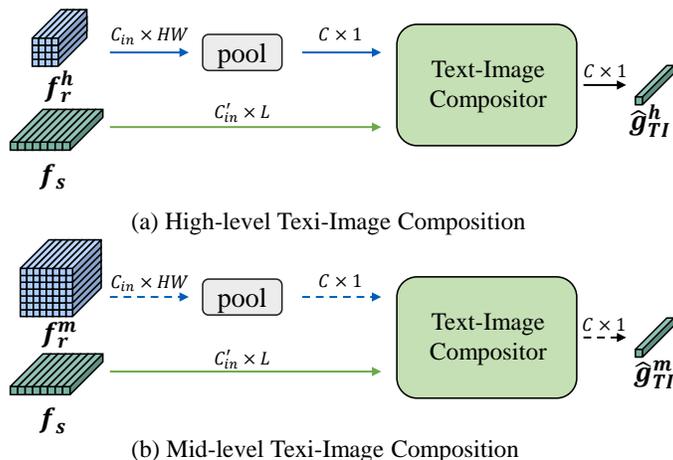


Figure 6: A brief illustration of two text-image compositors with the input shape. Please refer to Fig. 3 for the entire framework.

467 5.4. Analysis for hyperparameters

	R@1	R@10	R@50
Css-Net ($\alpha_{1-4} = 1, 1, 1, 1$)	20.04	56.44	80.87
Css-Net ($\alpha_{1-4} = 1, 0.5, 0.5, 0.5$)	20.13	56.81	81.32

Table 8: Ablation for hyperparameter α .

	R@1	R@10	R@50
Css-Net ($\lambda_{1-2} = 1, 1$)	19.95	56.64	80.55
Css-Net ($\lambda_{1-2} = 10, 1$)	20.13	56.81	81.32

Table 9: Ablation for Hyperparameter λ .

468 In this work, the hyperparameters α s and λ s are not handpicked, as we
 469 empirically find that they are not sensitive and do not affect the model perfor-
 470 mance significantly. We set α s to be 1:0.5:0.5:0.5 based on the observation
 471 that the high-level image-text compositor performs best among all composi-
 472 tors (Table 8) and we want this compositor to act like a leader in the group.
 473 Similarly, we use λ s = 1 for all compositors, as we have some preliminary
 474 experiments that show similar results with this setting (Table 9). To demon-
 475 strate this, we add some experimental results on the Shoes dataset, which
 476 is another challenging benchmark for composed image retrieval. The results
 477 show that our Css-Net achieves competitive performance with different val-
 478 ues of α s and λ s, indicating that our model is robust and stable to the choice
 479 of hyperparameters.

480

481 5.5. Effect of More Annotation Noise

482 In this work, we aim to relieve the issue of noisy annotations, which can
 483 compromise the entire training process. Further, we artificially increased the
 484 noise intensity during training by manually manipulating relevant captions,
 485 such as random deletion, random swap, and random insertion proposed in
 486 a NLP work (Wei and Zou, 2019). To be more specific, we conducted an
 487 experiment on the Shoes dataset for both the baseline and Css-Net. For each
 488 relative caption, there is a 50% probability of adding one of three types of
 489 noise: Each word in the sentence has a 50% probability of being deleted; half
 490 of the words in the sentence are replaced with synonyms; and new words
 491 are inserted into half of the word intervals. The performance of the newly
 492 developed baseline and Css-Net are shown in Table 10.

Method	R@1	R@10	R@50
Baseline (<i>w</i> noise)	16.29	50.14	75.91
Css-Net (<i>w</i> noise)	19.07	55.69	78.98
Baseline (<i>w/o</i> noise)	17.27	52.26	77.35
Css-Net (<i>w/o</i> noise)	20.13	56.81	81.32

Table 10: Effect of annotation noise (w/o refers to without; w refers to with).

493 6. Conclusion

494 We present a Consensus Network (Css-Net) for composed image retrieval.
495 Css-Net aims to relieve the inherent triplet ambiguity problem, which arises
496 when the dataset contains multiple false-negative candidates that match the
497 same query. This problem stems from annotators describing only simple
498 properties and frequently overlooking fine-grained details of the images. The
499 resulting noisy triplets significantly compromise the metric learning objective
500 and bias the single compositor. To this end, Css-Net employs a consensus
501 module with four compositors that possess distinct knowledge. As a group,
502 compositors learn mutually when training and infer collaboratively during
503 evaluation, effectively minimizing the negative effects caused by the triplet
504 ambiguity problem. Extensive experiments show that Css-Net has achieved
505 competitive recall performance on three widely-used benchmarks, without
506 substantially increasing the inference time.

507 References

- 508 Y. Guo, Z. Cheng, L. Nie, Y. Wang, J. Ma, M. Kankanhalli, Attentive long
509 short-term preference modeling for personalized product search, *ACM*
510 *Transactions on Information Systems* 37 (2019) 1–27.
- 511 R. Sharma, A. Vishvakarma, Retrieving similar e-commerce images using
512 deep learning, *arXiv:1901.03546* (2019).
- 513 Y. Guo, Z. Cheng, L. Nie, X.-S. Xu, M. Kankanhalli, Multi-modal preference
514 modeling for product search, in: *ACM Multimedia*, 2018.
- 515 H. Noh, A. Araujo, J. Sim, T. Weyand, B. Han, Large-scale image retrieval
516 with attentive deep local features, in: *ICCV*, 2017, pp. 3456–3465.
- 517 Z. Liu, P. Luo, S. Qiu, X. Wang, X. Tang, Deepfashion: Powering robust
518 clothes recognition and retrieval with rich annotations, in: *CVPR*, 2016,
519 pp. 1096–1104.
- 520 L. Liao, X. He, B. Zhao, C.-W. Ngo, T.-S. Chua, Interpretable multimodal
521 retrieval for fashion products, in: *ACM Multimedia*, 2018.
- 522 J. Deng, J. Guo, N. Xue, S. Zafeiriou, Arcface: Additive angular margin loss
523 for deep face recognition, in: *CVPR*, 2019, pp. 4690–4699.

- 524 X. Fan, W. Jiang, H. Luo, M. Fei, Spherereid: Deep hypersphere manifold
525 embedding for person re-identification, *Journal of Visual Communication*
526 *and Image Representation* 60 (2019) 51–58.
- 527 H. Sheng, Y. Zheng, W. Ke, D. Yu, X. Cheng, W. Lyu, Z. Xiong, Mining hard
528 samples globally and efficiently for person reidentification, *IEEE Internet*
529 *of Things Journal* 7 (2020) 9611–9622.
- 530 F. M. Hafner, A. Bhuyian, J. F. Kooij, E. Granger, Cross-modal distilla-
531 tion for rgb-depth person re-identification, *Computer Vision and Image*
532 *Understanding* 216 (2022) 103352.
- 533 L. Zhen, P. Hu, X. Wang, D. Peng, Deep supervised cross-modal retrieval,
534 in: *CVPR*, 2019, pp. 10394–10403.
- 535 Z. Zheng, L. Zheng, M. Garrett, Y. Yang, M. Xu, Y.-D. Shen, Dual-path
536 convolutional image-text embeddings with instance loss, *ACM Transac-*
537 *tions on Multimedia Computing, Communications, and Applications* 16
538 (2020) 1–23.
- 539 R. Guerrero, H. X. Pham, V. Pavlovic, Cross-modal retrieval and synthesis
540 (x-mrs): Closing the modality gap in shared subspace learning, in: *ACM*
541 *Multimedia*, 2021, pp. 3192–3201.
- 542 Z. Wang, Z. Gao, X. Xu, Y. Luo, Y. Yang, H. T. Shen, Point to rectangle
543 matching for image text retrieval, in: *ACM Multimedia*, 2022.
- 544 N. Vo, L. Jiang, C. Sun, K. Murphy, L.-J. Li, L. Fei-Fei, J. Hays, Composing
545 text and image for image retrieval an empirical odyssey, in: *CVPR*, 2019,
546 pp. 6439–6448.
- 547 Y. Chen, S. Gong, L. Bazzani, Image search with text feedback by visiolin-
548 guistic attention learning, in: *CVPR*, 2020, pp. 3001–3011.
- 549 S. Lee, D. Kim, B. Han, Cosmo: Content-style modulation for image retrieval
550 with text feedback, in: *CVPR*, 2021, pp. 802–812.
- 551 J. Kim, Y. Yu, H. Kim, G. Kim, Dual compositional learning in interactive
552 image retrieval, in: *AAAI*, volume 35, 2021, pp. 1771–1779.
- 553 H. Wen, X. Song, X. Yang, Y. Zhan, L. Nie, Comprehensive linguistic-visual
554 composition network for image retrieval, in: *SIGIR*, 2021, pp. 1369–1378.

- 555 A. Baldrati, M. Bertini, T. Uricchio, A. Del Bimbo, Conditioned and com-
556 posed image retrieval combining and partially fine-tuning clip-based fea-
557 tures, in: CVPR, 2022, pp. 4959–4968.
- 558 A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sas-
559 try, A. Askell, P. Mishkin, J. Clark, et al., Learning transferable visual
560 models from natural language supervision, in: ICML, 2021, pp. 8748–8763.
- 561 M. Wray, H. Doughty, D. Damen, On semantic similarity in video retrieval,
562 in: CVPR, 2021, pp. 3650–3660.
- 563 A. Falcon, G. Serra, O. Lanz, A feature-space multimodal data augmentation
564 technique for text-video retrieval, in: ACM Multimedia, 2022, pp. 4385–
565 4394.
- 566 V. B. Hinsz, Cognitive and consensus processes in group recognition memory
567 performance., *Journal of Personality and Social psychology* 59 (1990) 705.
- 568 S. Kullback, R. A. Leibler, On information and sufficiency, *The annals of*
569 *mathematical statistics* 22 (1951) 79–86.
- 570 T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, S. Belongie, Feature
571 pyramid networks for object detection, in: CVPR, 2017, pp. 2117–2125.
- 572 A. Miech, J.-B. Alayrac, I. Laptev, J. Sivic, A. Zisserman, Thinking fast and
573 slow: Efficient text-to-visual retrieval with transformers, in: CVPR, 2021,
574 pp. 9826–9836.
- 575 K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition,
576 in: CVPR, 2016, pp. 770–778.
- 577 Z. Zheng, L. Zheng, Y. Yang, A discriminatively learned cnn embedding
578 for person reidentification, *ACM transactions on multimedia computing,*
579 *communications, and applications* 14 (2017) 1–20.
- 580 Y. Sun, C. Cheng, Y. Zhang, C. Zhang, L. Zheng, Z. Wang, Y. Wei, Circle
581 loss: A unified perspective of pair similarity optimization, in: CVPR,
582 2020, pp. 6398–6407.
- 583 A. Wu, W.-S. Zheng, H.-X. Yu, S. Gong, J. Lai, Rgb-infrared cross-modality
584 person re-identification, in: ICCV, 2017, pp. 5390–5399.

- 585 P. Dai, R. Ji, H. Wang, Q. Wu, Y. Huang, Cross-modality person re-
586 identification with generative adversarial training., in: IJCAI, volume 1,
587 2018, p. 2.
- 588 J. Liu, Y. Sun, F. Zhu, H. Pei, Y. Yang, W. Li, Learning memory-augmented
589 unidirectional metrics for cross-modality person re-identification, in:
590 CVPR, 2022, pp. 19366–19375.
- 591 X. Qu, L. Liu, L. Zhu, L. Nie, H. Zhang, Source-free style-diversity
592 adversarial domain adaptation with privacy-preservation for person re-
593 identification, Knowledge-Based Systems 283 (2024) 111150.
- 594 R. Liu, Y. Zhao, S. Wei, L. Zheng, Y. Yang, Modality-invariant image-text
595 embedding for image-sentence matching, ACM Transactions on Multime-
596 dia Computing, Communications, and Applications 15 (2019) 1–19.
- 597 Q. Zhang, Z. Lei, Z. Zhang, S. Z. Li, Context-aware attention network for
598 image-text retrieval, in: CVPR, 2020, pp. 3536–3545.
- 599 Z. Liu, F. Chen, J. Xu, W. Pei, G. Lu, Image-text retrieval with cross-
600 modal semantic importance consistency, IEEE Transactions on Circuits
601 and Systems for Video Technology (2022).
- 602 L. Zhang, H. Wu, Q. Chen, Y. Deng, J. Siebert, Z. Li, Y. Han, D. Kong,
603 Z. Cao, Vldeformer: Vision–language decomposed transformer for fast
604 cross-modal retrieval, Knowledge-Based Systems 252 (2022) 109316.
- 605 Z. Li, C. Guo, X. Wang, H. Zhang, Y. Wang, Integrating listwise ranking
606 into pairwise-based image-text retrieval, Knowledge-Based Systems 287
607 (2024) 111431.
- 608 C. Deng, X. Xu, H. Wang, M. Yang, D. Tao, Progressive cross-modal seman-
609 tic network for zero-shot sketch-based image retrieval, IEEE Transactions
610 on Image Processing 29 (2020) 8892–8902.
- 611 H. Wang, C. Deng, T. Liu, D. Tao, Transferable coupled network for zero-
612 shot sketch-based image retrieval, IEEE Transactions on Pattern Analysis
613 and Machine Intelligence (2021) 1–1.
- 614 J. Li, Z. Ling, L. Niu, L. Zhang, Zero-shot sketch-based image retrieval with
615 structure-aware asymmetric disentanglement, Computer Vision and Image
616 Understanding 218 (2022) 103412.

- 617 S. Liang, W. Dai, Y. Cai, C. Xie, Sketch-based 3d shape retrieval via teacher-
618 student learning, *Computer Vision and Image Understanding* 239 (2024)
619 103903.
- 620 Y. Yang, M. Wang, W. Zhou, H. Li, Cross-modal joint prediction and align-
621 ment for composed query image retrieval, in: *ACM Multimedia*, 2021, pp.
622 3303–3311.
- 623 G. Zhang, S. Wei, H. Pang, Y. Zhao, Heterogeneous feature fusion and cross-
624 modal alignment for composed image retrieval, in: *ACM Multimedia*, 2021,
625 pp. 5353–5362.
- 626 C. Gu, J. Bu, Z. Zhang, Z. Yu, D. Ma, W. Wang, Image search with text
627 feedback by deep hierarchical attention mutual information maximization,
628 in: *ACM Multimedia*, New York, NY, USA, 2021.
- 629 Y. Zhao, Y. Song, Q. Jin, Progressive learning for image retrieval with
630 hybrid-modality queries, in: *SIGIR*, 2022, pp. 1012–1021.
- 631 X. Han, X. Zhu, L. Yu, L. Zhang, Y.-Z. Song, T. Xiang, Fame-vil: Multi-
632 tasking vision-language model for heterogeneous fashion tasks, in: *CVPR*,
633 2023, pp. 2669–2680.
- 634 M. Corbetta, G. L. Shulman, Control of goal-directed and stimulus-driven
635 attention in the brain, *Nature reviews neuroscience* 3 (2002) 201–215.
- 636 A. Blum, T. Mitchell, Combining labeled and unlabeled data with co-
637 training, in: *Proceedings of the eleventh annual conference on Compu-
638 tational learning theory*, 1998, pp. 92–100.
- 639 S. Qiao, W. Shen, Z. Zhang, B. Wang, A. Yuille, Deep co-training for semi-
640 supervised image recognition, in: *ECCV*, 2018, pp. 135–152.
- 641 J. Peng, G. Estrada, M. Pedersoli, C. Desrosiers, Deep co-training for semi-
642 supervised image segmentation, *Pattern Recognition* 107 (2020) 107269.
- 643 T. Hui, S. Liu, Z. Ding, S. Huang, G. Li, W. Wang, L. Liu, J. Han,
644 Language-aware spatial-temporal collaboration for referring video segmen-
645 tation, *IEEE Transactions on Pattern Analysis and Machine Intelligence*
646 (2023).

- 647 K. Saito, K. Watanabe, Y. Ushiku, T. Harada, Maximum classifier discrep-
648 ancy for unsupervised domain adaptation, in: CVPR, 2018, pp. 3723–3732.
- 649 Z. Zheng, Y. Yang, Unsupervised scene adaptation with memory regulariza-
650 tion in vivo, IJCAI (2019).
- 651 Y. Luo, L. Zheng, T. Guan, J. Yu, Y. Yang, Taking a closer look at domain
652 shift: Category-level adversaries for semantics consistent domain adapta-
653 tion, in: CVPR, 2019, pp. 2507–2516.
- 654 T. L. Berg, A. C. Berg, J. Shih, Automatic attribute discovery and charac-
655 terization from noisy web data, in: ECCV, 2010, pp. 663–676.
- 656 H. Wu, Y. Gao, X. Guo, Z. Al-Halah, S. Rennie, K. Grauman, R. Feris,
657 Fashion iq: A new dataset towards retrieving images by natural language
658 feedback, in: CVPR, 2021, pp. 11307–11317.
- 659 X. Han, Z. Wu, P. X. Huang, X. Zhang, M. Zhu, Y. Li, Y. Zhao, L. S. Davis,
660 Automatic spatially-aware fashion concept discovery, in: ICCV, 2017, pp.
661 1463–1471.
- 662 J.-H. Kim, S.-W. Lee, D. Kwak, M.-O. Heo, J. Kim, J.-W. Ha, B.-T. Zhang,
663 Multimodal residual learning for visual qa, in: NeurIPS, volume 29, 2016.
- 664 E. Perez, F. Strub, H. De Vries, V. Dumoulin, A. Courville, Film: Visual
665 reasoning with a general conditioning layer, in: AAAI, volume 32, 2018.
- 666 G. Delmas, R. S. Rezende, G. Csurka, D. Larlus, Artemis: Attention-based
667 retrieval with text-explicit matching and implicit similarity, in: ICLR,
668 2022.
- 669 Y. Chen, Z. Zheng, W. Ji, L. Qu, T.-S. Chua, Composed image re-
670 trieval with text feedback via multi-grained uncertainty regularization,
671 arXiv:2211.07394 (2022).
- 672 A. Hermans, L. Beyer, B. Leibe, In defense of the triplet loss for person
673 re-identification, arXiv:1703.07737 (2017).
- 674 X. Wang, H. Zhang, W. Huang, M. R. Scott, Cross-batch memory for em-
675 bedding learning, in: CVPR, 2020, pp. 6388–6397.

- 676 T. Chen, S. Kornblith, M. Norouzi, G. Hinton, A simple framework for
677 contrastive learning of visual representations, in: ICML, 2020, pp. 1597–
678 1607.
- 679 K. He, H. Fan, Y. Wu, S. Xie, R. Girshick, Momentum contrast for unsu-
680 pervised visual representation learning, in: CVPR, 2020, pp. 9729–9738.
- 681 S. Hochreiter, J. Schmidhuber, Long short-term memory, *Neural computa-*
682 *tion* 9 (1997) 1735–1780.
- 683 F. Zhang, M. Yan, J. Zhang, C. Xu, Comprehensive relationship reasoning
684 for composed query based image retrieval, in: ACM Multimedia, 2022.
- 685 R. Girshick, Fast r-cnn, in: ICCV, 2015, pp. 1440–1448.
- 686 A. Graves, Long short-term memory, *Supervised sequence labelling with*
687 *recurrent neural networks* (2012) 37–45.
- 688 Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis,
689 L. Zettlemoyer, V. Stoyanov, Roberta: A robustly optimized bert pre-
690 training approach, arXiv:1907.11692 (2019).
- 691 D. P. Kingma, J. Ba, Adam: A method for stochastic optimization,
692 arXiv:1412.6980 (2014).
- 693 J. Wei, K. Zou, Eda: Easy data augmentation techniques for boosting perfor-
694 mance on text classification tasks, in: EMNLP-IJCNLP, 2019, pp. 6382–
695 6388.