

# L<sup>A</sup>T<sub>E</sub>X Author Guidelines for CVPR Proceedings

Anonymous CVPR submission

Paper ID \*\*\*\*

## Abstract

pass

## 1. Introduction

The importance of our research area

Some progress in sketch based image retrieval

Difficulty in Zero-shot setup and some possible solutions

Our proposed methods and their advantages

itemize our contributions in this paper

## 2. Related Work

### 2.1. Sketch-based image retrieval

### 2.2. Zero-Shot Learning

### 2.3. Cross-modal domain translation

## 3. Methodology

There will be five parts in this section. Sec. 3.1 defines the our targeted problem and briefly introduce our framework. Sec.3.2 introduce the feature extractor. Sec. 3.3 introduce the image cVAE-GAN. Sec. 3.4 introduce the sketch cVAE-GAN. Sec. 3.5 introduce the semantic preservation module. Sec. 3.6 introduce the design of loss functions during training procedure.

### 3.1. Problem Definition

In this paper, we focus on solving the problem of hand-free sketch-based image retrieval under zero-shot setup, where only the sketches and images from seen class are used during training stage. Our proposed framework is expected to use the sketches to retrieve the images, the categories of which have never appeared during training.

We first provide a definition of the SBIR in zero-shot setting. Given a dataset  $S = \{(x_i^{img}, x_i^{ske}, x_i^{sem}, y_i) | y_i \in \mathcal{Y}\}$ , where  $x_i^{img}$ ,  $x_i^{ske}$ ,  $x_i^{sem}$  and  $y_i$  are corresponding to the image, sketch, semantic representation and class label. Following the zero-shot setting in [2], we split all classes  $\mathcal{Y}$  into  $\mathcal{Y}_{train}$  and  $\mathcal{Y}_{test}$  according to whether the label exists in ImageNet[1], where no overlap exists between two label set, i.e.  $\mathcal{Y}_{train} \cap \mathcal{Y}_{test} = \emptyset$ . Based on the partition of label set  $\mathcal{Y}$ , we split dataset into  $S_{train}$  and  $S_{test}$ . Our model need to disentangle structure representations of image using data in  $S_{train}$ . During test, given  $x^{ske}$  from  $S_{test}$ , our model need to retrieve several images from test images candidate.

Our goal is to learn a two-way map between image feature domain to sketch feature domain. To this end, we propose a new deep network (shown in Figure ??), which contains two structure encoders  $\{E_s^{img}, E_s^{ske}\}$ , one appearance encoder  $E_a^{img}$ , two feature decoder  $\{G^{img}, G^{ske}\}$ , a semantic decoder and two domain discriminators  $\{D^{img}, D^{ske}\}$ . Note that, the overall model can be regarded as two cVAE-GANs working parallel, which target to reconstruct sketch features from image and reconstruct image features from both sketches and images. To better capture the semantics information inside the sketches and images, we also add a semantic decoder to preserve semantics information while reconstructing the image features.

### 3.2. Feature Extractor

Considering the abstractness and visual sparsity of sketch, it is challenging extract feature from sketch. To alleviate this issue, multi-channel and multi-scale model was proposed to extract more saint features [3]. Motivated by the visualization in [4], where different layers capture visual features at different levels, we follow [] and build our feature extractor using a multi-layer feature fusion network enrich feature representation capacity without adding any additional parameters.

108	<b>3.3. Image cVAE-GAN</b>	162
109	<b>3.4. Sketch cVAE-GAN</b>	163
110	<b>3.5. Semantics Preservation</b>	164
111	<b>3.6. Loss Function</b>	165
112	<b>4. Experiment</b>	166
113	<b>4.1. Experiment Setup</b>	167
114	<b>4.1.1 Dataset</b>	168
115	<b>4.1.2 Implementation Details</b>	169
116	<b>4.2. Comparison</b>	170
117	<b>4.3. Ablation Study</b>	171
118	<b>4.4. Case study</b>	172
119	<b>5. Conclusion</b>	173
120	<b>6. To Discuss</b>	174
121	<ul style="list-style-type: none"><li>• Whether to generator the whole image/sketch.</li></ul>	175
122	<ul style="list-style-type: none"><li>• If the poses between the image and the sketch are different, can the model learn the sketch information between image and sketch.</li></ul>	176
123	<ul style="list-style-type: none"><li>• Where to add the semantics information to further supervise the model's training.</li></ul>	177
124	<b>References</b>	178
125	[1] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In <i>2009 IEEE conference on computer vision and pattern recognition</i> , pages 248–255. IEEE, 2009. 1	179
126	[2] Sasi Kiran Yelamarthi, Shiva Krishna Reddy, Ashish Mishra, and Anurag Mittal. A zero-shot framework for sketch based image retrieval. In <i>European Conference on Computer Vision</i> , pages 316–333. Springer, 2018. 1	180
127	[3] Qian Yu, Yongxin Yang, Feng Liu, Yi-Zhe Song, Tao Xiang, and Timothy M Hospedales. Sketch-a-net: A deep neural network that beats humans. <i>International journal of computer vision</i> , 122(3):411–425, 2017. 1	181
128	[4] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In <i>European conference on computer vision</i> , pages 818–833. Springer, 2014. 1	182
129		183
130		184
131		185
132		186
133		187
134		188
135		189
136		190
137		191
138		192
139		193
140		194
141		195
142		196
143		197
144		198
145		199
146		200
147		201
148		202
149		203
150		204
151		205
152		206
153		207
154		208
155		209
156		210
157		211
158		212
159		213
160		214
161		215