

L^AT_EX 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 parallel cVAE-GANs. Sec. 3.4 introduce the semantic preservation module. Sec. 3.5 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 [?], we split all classes \mathcal{Y} into \mathcal{Y}_{train} and \mathcal{Y}_{test} according to whether the label exists in ImageNet[?], 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 feature extractor $F(\cdot)$, two structure encoders $\{E_s^{img}(\cdot), E_s^{ske}(\cdot)\}$, one appearance encoder $E_a^{img}(\cdot)$, two feature decoder $\{G^{img}(\cdot, \cdot), G^{ske}(\cdot, \cdot)\}$, a semantic decoder $S(\cdot)$ and two domain discriminators $\{D^{img}(\cdot, \cdot), D^{ske}(\cdot, \cdot)\}$. Note that by combining the encoders, decoders and discriminators, our 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 preserve the semantics information within 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 [?]. Motivated by the visualization in [?], where different layers in the backbone model capture visual features at different levels, we build our feature extractor by combining different layers' features together to enrich feature representation capacity without adding any additional parameters.

In detail, we first use a pre-trained backbone model on ImageNet [?] to process each sketch and image, where we

use VGG-16 as our backbone model in this paper. Suppose $f_i \in \mathbb{R}^{C_i \times H_i \times W_i}$ is the output of the i -th convolution module and $f_{fc} \in \mathbb{R}^N$ is the feature of the last fully connected layer, the extracted feature f^* can be formulated as:

$$f^* = F(x_i) = [f_{fc}, GAP(f_5), GAP(f_4), GAP(f_3)] \quad (1)$$

where $GAP(\cdot)$ means global spatial average pooling and $[\cdot, \cdot]$ means the concatenation operation between vectors. Following the instruction in [?], a Principal Component Analysis algorithm is adopted to reduce the dimension of the extracted feature, which is helpful to not only improve the computational efficiency but remove some redundant information.

3.3. Parallel cVAE-GAN

Generative approaches have shown their power in ZS-SBIR [?, ?]. But both of these two papers, only encode the image in one encoder, which makes the appearance and the structure features fused together. When training the VAE based model, the decoder may ignore the condition information in sketch feature. Motivated by [?], we can regard image and sketch as two domains, where image domain have richer information than sketch domain. To this end, we use three encoders to encode image and sketch features respectively. Given

3.4. Semantics Preservation

3.5. Loss Function

4. Experiment

4.1. Experiment Setup

4.1.1 Dataset

4.1.2 Implementation Details

4.2. Comparison

4.3. Ablation Study

4.4. Case study

5. Conclusion

6. To Discuss

- Whether to generator the whole image/sketch.
- If the poses between the image and the sketch are different, can the model learn the sketch information between image and sketch.
- Where to add the semantics information to further supervise the model's training.