

汇报人：江篱

文献阅读

Sketch Less for More: On-the-Fly Fine-Grained Sketch Based Image Retrieval

文献信息：Bhunia, Ayan Kumar, et al. "Sketch Less for More: On-the-Fly Fine-Grained Sketch Based Image Retrieval." *arXiv preprint arXiv:2002.10310* (2020).

问题：手绘草图检索图像时，绘制草图需要时间，而且大多数人并不能绘制一幅完整且形状相似的草图。

解决办法：我们重新制定了FG-SBIR（Fine-grained sketch-based image retrieval）框架，目标是以尽可能少的笔画检索目标照片。一旦用户开始绘图，就开始检索。具体实现方式是设计一个基于强化学习的跨模态检索框架，优化的时候用完整的素描图像检索出的ground truth图像的排序进行优化。这篇论文还提出一种新的奖励方案，该方案规避了与无关笔画相关的问题，从而在检索过程中为模型提供更一致的图像排名。

contributions:

1. 提出了on-the-fly FG-SBIR框架，使用强化学习利用不完整的草图检索图像。
2. 新的损失函数。

概述：下图展示了本文提出方法与当前最好的FG-SBIR 模型方法（损失函数为triplet loss）的结果比较。要正确结果出现在前十个答案中，本文方法在用户完成了30%的草图时便可得到这个结果，我们选择的基准模型要在用户完成80%的草图才能得到该结果。

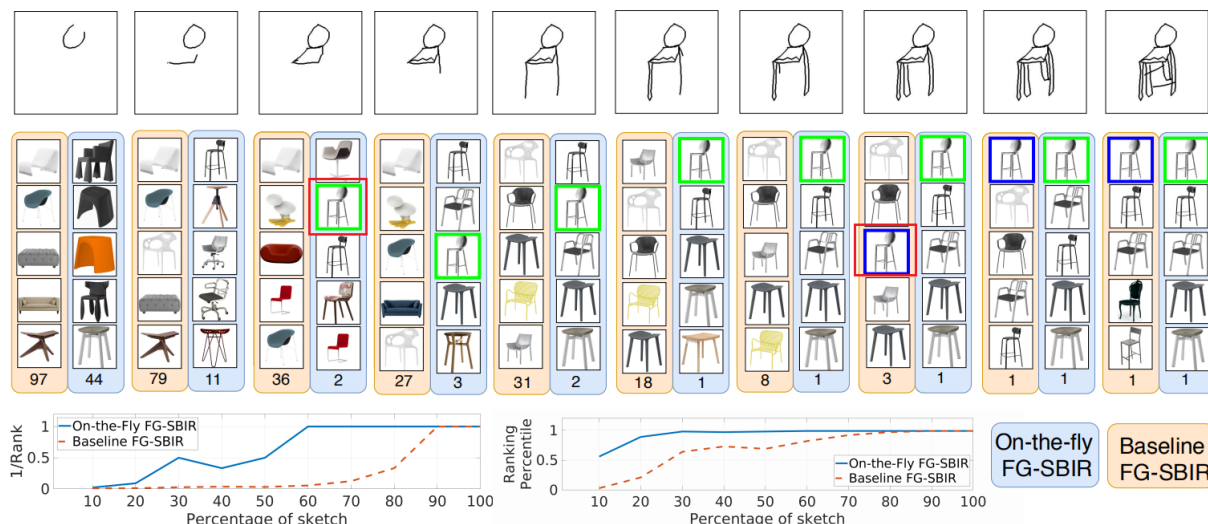


Figure 2. Illustration of proposed *on-the-fly* framework's efficacy over a baseline FG-SBIR method [41, 49] trained with completed sketches only. For this particular example, our method needs only 30% of the complete sketch to include the true match in the top-10 rank list, compared to 80% for the baseline. Top-5 photo images retrieved by either framework are shown here, in progressive sketch-rendering steps of 10%. The number at the bottom denotes the paired (true match) photo's rank at every stage.

Method：我们的目的是根据草图在图像库中检索出相似图像，以下用公式来阐述我们的目的：

1. $G = \{x_i\}_{i=1}^M$ $\xrightarrow{\text{特征提取}}$ $\hat{G} = \{F(x_i)\}_{i=1}^M$ G 图像库 特征向量
2. query sketch S $\xrightarrow{\text{Ret}_q(F(S), \hat{G})}$ 根据草图S在G中检索出q个最相似的图像
3. $S \in \{p^1, p^2, p^3, \dots, p^N\}$ $\xrightarrow{\text{Ret}_q(F(\phi(S^k)), \hat{G})}$ 多个未完成的草图 我们的任务

模型：下图为传统模型与本文方法的对比。

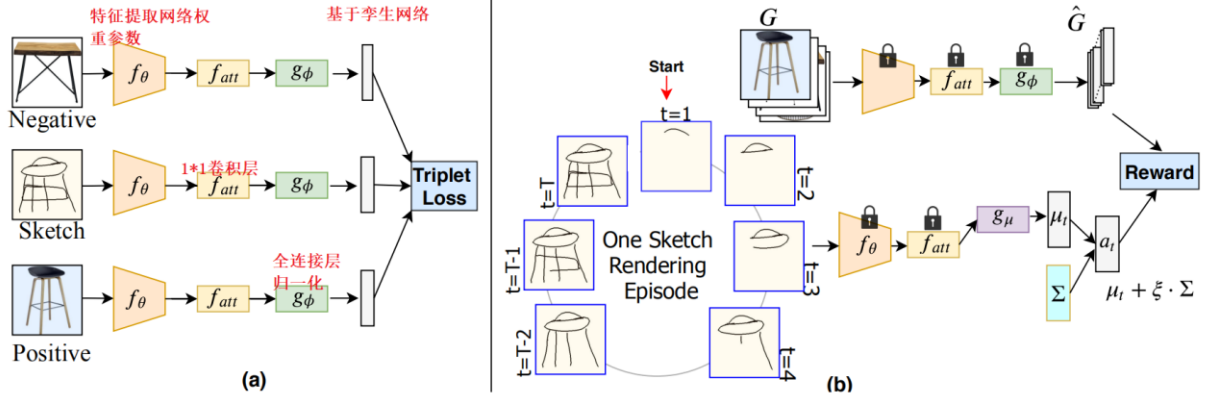


Figure 3. (a) A conventional FG-SBIR framework trained using triplet loss. (b) Our proposed reinforcement learning based framework that takes into account a complete sketch rendering episode. Key locks signifies particular weights are fixed during RL training.

传统模型：

我们将参数 f_θ 、 f_{att} 和 g_ϕ 所代表的过程打包为一个整体编码函数 F 。网络的输入为 $\{a, p, n\}$ ， a 表示手绘草图， p 表示与草图相似的图像（positive image）， n 表示不相似的图像（negative image）。

草图与positive image的距离为 $\beta^+ = \|F(a) - F(p)\|_2$ 。

草图与negative image的距离为 $\beta^- = \|F(a) - F(n)\|_2$ 。

损失函数为 $\max\{0, \mu + \beta^+ - \beta^-\}$ ， μ 为超参数。

本文方法：

本文方法如上图（b）部分所示，基准为相似图像 G 的特征表示 \hat{G} ，输入为草图的整个绘制过程 $S \in \{p_1, p_2, p_3, \dots, p_T\}$ ，对于每一个输入 p_t ，输出为一个特征向量 a_t 。损失函数即为 a_t 与 \hat{G} 的距离。

遵循典型的强化学习表示法，我们将我们的策略定义为 $\pi_\theta(a|s)$ 。

$$\pi_{\Theta}(a_t|s_t) = \sqrt{\frac{1}{(2\pi)^D |\Sigma|}} \times \exp \left\{ -\frac{1}{2} (a_t - \mu_t)^{\top} \Sigma^{-1} (a_t - \mu_t) \right\}, \quad (1)$$

where the mean $\mu_t = g_{\mu}(s'_t) \in \mathbb{R}^D$, and s'_t is obtained via a pre-trained f_{θ} and f_{att} that take state $s_t = \mathcal{O}(p_t)$ as its input. Meanwhile, Σ is a standalone trainable diagonal covariance matrix. We sample action $a_t = \mu_t + \xi \cdot \Sigma$, where $\xi \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $a_t \in \mathbb{R}^D$.

就损失函数来说分为Local Reward和Global Reward。

Local Reward: 损失函数如下，当rank排名越靠后， R_t^{Local} 值越小，在训练的过程中，我们最大化该函数。

$$R_t^{Local} = \frac{1}{rank_t}$$

Global Reward: 草图完成的过程中有很多步，假设某一步用 L_t 来表示，考虑到早期草图随机性更高，那么在草图完成的过程中，后期两步的距离要小于前期两步的距离。用公式来表示为：

$$R_t^{Global} = -\max(0, \tau(L_t, L_{t+1}) - \tau(L_{t-1}, L_t))$$