

StyleMeUp: Towards Style-Agnostic Sketch-Based Image Retrieval

Aneeshan Sain^{1,2} Ayan Kumar Bhunia¹ Yongxin Yang^{1,2} Tao Xiang^{1,2} Yi-Zhe Song^{1,2}

¹SketchX, CVSSP, University of Surrey, United Kingdom.

²iFlyTek-Surrey Joint Research Centre on Artificial Intelligence.

{a.sain, a.bhunia, yongxin.yang, t.xiang, y.song}@surrey.ac.uk

Abstract

Sketch-based image retrieval (SBIR) is a cross-modal matching problem which is typically solved by learning a joint embedding space where the semantic content shared between photo and sketch modalities are preserved. However, a fundamental challenge in SBIR has been largely ignored so far; that is, sketches are drawn by humans and considerable style variations exist amongst different users. An effective SBIR model needs to explicitly account for this style diversity, crucially, to generalise to unseen user styles. To this end, a novel style-agnostic SBIR model is proposed. Different from existing models, a cross-modal variational autoencoder (VAE) is employed to explicitly disentangle each sketch into a semantic content part shared with the corresponding photo, and a style part unique to the sketcher. Importantly, to make our model dynamically adaptable to any unseen user styles, we propose to meta-train our cross-modal VAE by adding two style-adaptive components: a set of feature transformation layers to its encoder and a regulariser to the disentangled semantic content latent code. With this meta-learning framework, our model can not only disentangle the cross-modal shared semantic content for SBIR, but can adapt the disentanglement to any unseen user style as well, making the SBIR model truly style-agnostic. Extensive experiments show that our style-agnostic model yields state-of-the-art performance for both category-level and instance-level SBIR.

1. Introduction

Sketch as an input modality has been proven to be a worthy complement to text for photo image retrieval [8, 10, 4]. Its precision in visual description is particularly useful for fine-grained retrieval, where the goal is to find a specific object instance rather than category [44, 49, 3]. Research has flourished in recent years, where the main focus has been on addressing the sketch-photo domain gap [33, 16, 44] and data scarcity [5, 3, 10, 35, 12]. Thanks to these combined efforts, reported performances have already shown promise for practical adaptation.

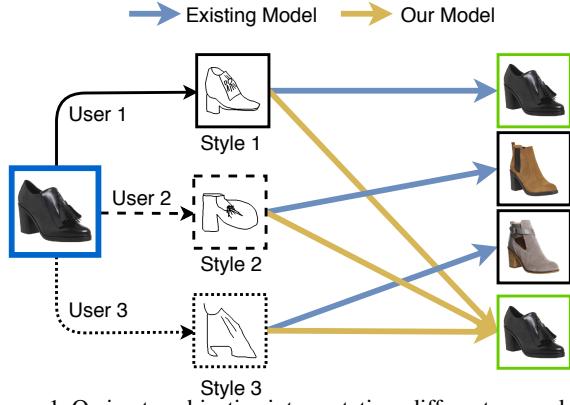


Figure 1. Owing to subjective interpretation, different users sketch the same object instance (a shoe here) very differently. Without considering this style diversity, an existing SBIR model yields completely different results for different sketches. With our style-agnostic model, the same intended object is retrieved.

However, there is an important issue that has largely been ignored so far and has impeded the effectiveness of existing SBIR models – sketches are drawn by humans and there exists considerable style variations amongst users (Fig. 1). This is a result of subjective interpretation and different drawing skills of different users. Consequently, even with the same object instance as reference, sketches of different users can look drastically different as shown in the example in Fig. 1. Existing SBIR models [30, 8, 47, 35, 44] focus primarily on bridging the gap between the photo and sketch modalities. This is typically achieved by learning a joint embedding space where only the common semantic content part of a matching photo-sketch pair are preserved for matching. However, the large style variations of different users mean that the shared common semantic content can also vary (e.g., Fig. 1 shows that different users may choose to depict different characteristics of the same shoe). Crucially, it can vary in an unpredictable way – a commercial SBIR model will be used mostly by users whose sketches have never been used for model training. These models are thus poorly equipped to cope with this style diversity and unable to generalise to new user styles.

In this paper, a novel style-agnostic SBIR framework is

proposed which explicitly accounts for the style diversity and importantly can adapt dynamically to any unseen user styles without any model retraining. Different from existing SBIR models which focus solely on the shared semantic content between the photo and sketch modality and discard the modal-specific parts, we argue that in order to effectively deal with the style variations unique to the sketch modality, a disentanglement model is needed. With such a model, the user style can be modelled explicitly, making way for better generalisation.

The core of our style-agnostic SBIR framework is thus a disentanglement model, that takes a sketch or photo image as input and decomposes its content into a cross-modal shared semantic part to be used for retrieval, and a modal-specific part – in case of sketch, it corresponds to the user’s drawing style. Disentangling sketching styles is however a challenging task. Existing style disentanglement methods usually cater to problems where the style information carry less variance (e.g., schools of art, building styles, etc) and hence is comparatively easier to separate [24]. For sketches however, we are faced with much larger *variability* where each user has a *unique* style that can manifest itself in different ways for different object instances. The disentanglement should thus be able to dynamically adapt to new user styles and new object categories for better generalisation. To this end, we propose a novel disentanglement model with meta-learning, that generalises to unseen user styles.

Concretely, we employ a cross-modal translation variational autoencoder (VAE) framework [23] to project a sketch/photo into its modal-invariant semantic part, and modal-specific part. The VAE is used for both sketch reconstruction and translation as well as sketch-to-photo translation to exploit the shared semantic content across both modalities and styles. To make the disentanglement dynamically adaptable and generalisable, a popular gradient-based meta-learning model namely model-agnostic meta-learning (MAML) [14] is adopted. Designed for few-shot learning, the original MAML cannot be directly applied. We thus introduce two new components as shown in Fig. 2: (i) A set of feature transformation layers sitting between the VAE encoder layers for the encoder adaptation, and (ii) a regulariser designed to adapt the disentangled modal-invariant semantic content part of the latent code produced by the encoder. Both component’s parameters are meta-learned using MAML for fast adaptation to new style/categories sampled during episodic training of MAML. Once trained, these two components are responsible for adaptation to new/unseen user styles and object categories/instances, therefore achieving style-agnostic SBIR.

Our contributions are as follows: (a) For the first time, we propose the concept of style-agnostic SBIR to deal with a largely neglected user style diversity issue in SBIR. (b) We introduce a novel style-agnostic SBIR framework

based on disentangling a photo/sketch image into a modal-invariant semantic content part suitable for SBIR and a model-specific part that needs to be explicitly modelled in order to minimise its detrimental effects on retrieval. (c) To make the disentanglement generalisable to unseen user styles and object categories/instances, feature transformation layers and latent modal-invariant code regulariser are introduced to a VAE, both of which are meta-learned using a MAML framework for style/category/instance adaptation. (d) Extensive experiments show that state-of-the-art performances can be achieved as a direct result of the style-agnostic design.

2. Related Works

Category-level SBIR: Category-level SBIR tasks accept a sketch-query with an aim to retrieve photos of the same category [45, 9]. Early approaches deploy handcrafted descriptors [52] such as SIFT [31], Gradient Field HOG [19], Histogram of Edge Local Orientations [42] or Learned Key Shapes [43], for constructing local [19] or global [40] joint photo-sketch representations. Most recent approaches are based on deep learning [30, 8]. They typically employ Siamese-like neural networks with ranking losses, like triplet loss [60] to learn a joint embedding space for both the sketch and photo modalities. Contemporary research directions also include zero-shot SBIR [59, 10] where a model aims to generalise across disjoint training and testing classes [12], alleviating annotation costs. Sketch-photo hashing [30, 63] on the other hand embeds to binary hash-codes instead of continuous vectors for computational ease.

Fine-grained SBIR: As opposed to category-level SBIR, fine-grained SBIR (FG-SBIR) [47, 35, 44] is directed towards instance-level sketch-photo matching. Starting with deformable-part models [27], various deep approaches have surfaced with the advent of new FG-SBIR datasets [49, 60, 20]. Yu *et al.* [60] introduced a deep triplet-ranking model that learnt a joint sketch-photo embedding space. This was further enhanced via attention based techniques with advanced higher order retrieval loss [49], hybrid generative-discriminative cross-domain image generation [36], textual tags [48] and employing mixed modal jigsaw solving for a better pre-training strategy [37]. While Sain *et al.* [44] explored cross-modal hierarchical co-attention amongst sketch-photo regions, Bhunia *et al.* [5] employed reinforcement learning in an early retrieval scenario. These models focus on learning a joint embedding space where only the modal-invariant shared semantic content of a matching photo-sketch pair is preserved for both modalities. However, without explicitly modelling the modal-specific parts, particularly for sketches the user styles, these models cannot generalise well to unseen objects and user styles.

Disentangled representation learning: Learning a disentangled representation would require modelling dis-

tinct informative factors in the variations of data [11]. Starting from generic frameworks like combining auto-encoders with adversarial training [32], this disentanglement paradigm has been successfully applied to recognition [55, 38], image-to-image translation [58] and image-editing [22, 56] tasks. While, InfoGAN [7] optimises mutual knowledge between latent variables, β -VAE [17] balances a hyperparameter β to learn independent data generative factors for disentanglement in an unsupervised setting. These methods however lack interpretability, with the relevance of each learned factor being uncontrollable. A few recent works include joint disentanglement and adaptation module trained in a cross-domain cycle-consistency paradigm [64] or multi-scale spatial-temporal maps with a cross-verified disentangling strategy [34] or adversarial parameter estimation [39]. None of such methods however has worked towards disentangling features for sketches to assist in SBIR. Furthermore, none has the ability to adapt the disentanglement dynamically for new user styles and object instances, which our meta-learning based cross-modal disentanglement VAE is designed for.

Meta-Learning: Meta-learning aims to acquire transferable knowledge from different sample training-tasks to help adapt to unseen tasks with only a few training samples [18]. Most existing meta-learning methods [54, 46, 6] are designed for few-shot image classification and thus are suitable for our problem of meta-learning of a generalisable cross-modal disentanglement model. The popular gradient/optimisation based meta-learning method MAML [14], however is general enough to be adapted to our problem. MAML trains a base model on a set of source tasks to learn good initialisation parameters that adapts quickly to new tasks during training, over just a few gradient descent updates. Since its first introduction, various modifications have been proposed including Meta-SGD [28], MAML++ [1], latent embedding optimization (LEO) [41] and uncertainty-induced MAML for continual learning [15]. Among them, those designed for domain adaptation [53] or domain generalization [26, 2] are the most relevant to our work. Different from them, our model uniquely addresses the cross-modal SBIR problem via meta-learning a disentanglement VAE for *generalising* better onto unseen user styles and object instances.

3. Methodology

Overview: We aim to devise a SBIR framework that learns to model the diversity in sketching-styles corresponding to the same object category (for category-level SBIR [9]) or instance (for FG-SBIR [44]). To this end, we design a style agnostic disentanglement model which decomposes the content of a photo/sketch image into a modal-invariant semantic part suitable for cross-modal matching, and a modal-specific part which is a distractor to SBIR but

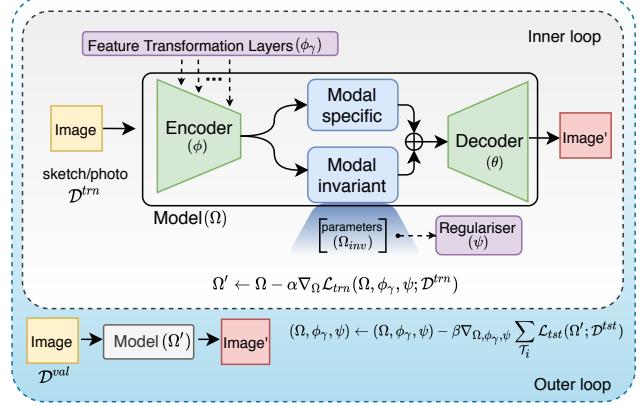


Figure 2. Our core model is a VAE framework that disentangles the modal variant and invariant semantics in a sketch in a cross-modal translation setting. While a regulariser network regularises parameters of the invariant component (Ω_{inv}), feature transformation (FT) layers aid in style-agnostic encoding following a meta-learning paradigm.

needs to be modelled explicitly to assist in the disentanglement. The disentanglement model is a cross-modal VAE that learns to embed a photo/sketch image to reconstruct them either in its original modality or to the other modality.

Formally, we are given a set of $C = \{C_1, C_2, \dots, C_M\}$ categories ($M \geq 1$) where every category C_i has $d^i = \{d_1^i, d_2^i, \dots, d_{N_i}^i\}$ ($N_i \geq 1$) data-points. Every data-point $d_j^i|_{j=1}^{N_i}$ corresponds to a sketch(s)-photo(p) pair i.e. $\{s_j^i, p_j^i\}$. For FG-SBIR, every photo instance is treated as a category with multiple sketch-styles paired with the same photo (p^i), i.e. $\forall d_j^i, j \in [1, N^i] \ p_j = p^i$. Feeding each data point to the encoder of the VAE, a latent code is obtained for both the photo and sketch. Our model (Fig. 2) aims to disentangle the latent code into modal-invariant and modal-specific components. The former corresponds to the semantic content of the object and thus should be used for cross-modal matching. It is subjected to a triplet loss so as to minimise the distance of from a sketch sample (s) to its matching photo sample ($p+$), while increasing that to an unmatched one ($p-$). Such a model is trained in a meta-learning framework for better generalisation. Once trained, during inference it uses the learnt encoder and the modal-invariant component to produce a style-agnostic embedding function $\mathcal{F}(\cdot) : \mathbb{R}^{H \times W \times 3} \rightarrow \mathbb{R}^{d_h}$ to map a rasterised sketch/photo having height H and width W to a \mathbb{R}^{d_h} feature for matching.

3.1. Disentanglement by Cross-modal Translation

Our disentanglement model is built upon a VAE framework [23] for both intra-modal reconstruction and cross-modal translation. The original VAE model produces a latent representation by optimising the variational lower bound on log-likelihood of the data

$$\log \mathbf{p}(x) \geq \mathbb{E}_{z \sim \mathbf{q}(z|x)} [\log \mathbf{p}(x|z)] - D_{\text{KL}}(\mathbf{q}(z|x) || p(z)), \quad (1)$$

where $D_{\text{KL}}(\cdot, \cdot)$ is the Kullback-Leibler (KL) divergence, and the conditional probability distributions $\mathbf{q}(z|x)$, $\mathbf{p}(x|z)$ refer to the encoder and decoder respectively, both parameterised by neural networks. The distribution $\mathbf{p}(z)$ is the prior on the latent space, modeled as $\mathcal{N}(z|\mathbf{0}, \mathbf{I})$. The encoder returns mean μ and variance σ^2 of a normal distribution, such that $z \sim \mathcal{N}(\mu, \sigma^2)$. Unlike this formulation which takes a single distribution into account, cross-modal training requires considering at least two modalities. Following [50] we extend Eq. 1 to the multi-modality case as:

$$\log \mathbf{p}(x_o) \geq \mathbb{E}_{z \sim \mathbf{q}(z|x_i)} [\log \mathbf{p}(x_o|z)] - D_{\text{KL}}(\mathbf{q}(z|x_i) || p(z)), \quad (2)$$

where x_i, x_o represents samples from input and output modalities respectively. It shows that the input and output samples can be decoupled via a joint embedding space shared by both sketch and photo modalities. This shared space thus allows for both same as well as cross-modality translations. Cross-modal training importantly creates a manifold which places objects to a high-dimensional space enriched with sketch-photo semantic relevance.

So far, we have described a VAE for cross-modal as well as intra-modal image translation. Our objective is however to disentangle the image content into modal-invariant and model-specific parts. Such a disentanglement takes place in the latent space produced by the encoder. More specifically, our CNN encoder Enc_ϕ projects an input I into two parts: a modal-invariant component (z_{inv}), and parameters mean (μ) and variance (σ) for a variable (modal-specific) component (z_{var}). Essentially the variable component is modelled via an independent unit Gaussian distribution as $z_{var} = \mu + \sigma \odot \mathcal{N}(0, 1)$ where $z \in \mathbb{R}^{d_h}$. Combining both components we thus obtain our final latent variable $z_f = z_{var} \oplus z_{inv}$, where \oplus represents element-wise summation. z_f is then fed to the decoder for reconstruction as $\hat{I} = Dec_\theta(z_f)$. Such a VAE model is trained by optimising the sum of reconstruction (\mathcal{L}_{rec}) and KL divergence (\mathcal{L}_{KL}) losses via gradient descent, with:

$$\begin{aligned} \mathcal{L}_{\text{rec}}(\phi, \theta) &= -\mathbb{E}_{\mathbf{q}(z_f|I)} [\log \mathbf{p}(I|z_f)], \\ \mathcal{L}_{\text{KL}} &= D_{\text{KL}}[\mathbf{q}_\phi(z_f|I) || \mathbf{p}(z_f)], \end{aligned} \quad (3)$$

where the prior over latent variables is a centered isotropic multivariate Gaussian, $\mathbf{p}_\theta(z_f) = \mathcal{N}(z_f; \mathbf{0}, I)$. In practice however, we simplify $\mathcal{L}_{\text{rec}} = \|\hat{I} - I\|_2$.

Besides sketch-photo translation we perform cross-style translation between two sketches of the same object, which ensures modelling the style diversity. Concretely, given two sketches s_j and s_k we model the latent feature of s_j as $z_f^{s_j*} = z_{inv} \oplus z_{var}^{s_k}$, and then we reconstruct it by $\hat{s}_j^* = Dec_\theta(z_f^{s_j*})$, where s_k is another randomly chosen style of the same object ($j, k \in [1, N_i]; j \neq k$). Accordingly we obtain the sum of all reconstruction losses as \mathcal{L}_{rec}

and sum their corresponding KL-divergence losses as \mathcal{L}_{KL} . To instil discriminative knowledge, we train the invariant component with a triplet-loss ($\mathcal{L}_{\text{Tri}}^{z_{inv}}$) objective [60] where the distance of z_{inv} extracted from a sketch (denoted as s), is reduced from that of its matching photo ($p+$), and increased from that of a non-matching one ($p-$). Furthermore, we resort towards discriminative sample generation by imposing a similar triplet objective on the synthesised embedding features z_f . We thus have:

$$\begin{aligned} \mathcal{L}_{\text{Tri}}^{z_{inv}} &= \max\{0, m^{z_{inv}} + \delta(z_{inv}^s, z_{inv}^{p+}) - \delta(z_{inv}^s, z_{inv}^{p-})\}, \\ \mathcal{L}_{\text{Tri}}^{z_f} &= \max\{0, m^{z_f} + \delta(z_f^s, z_f^{p+}) - \delta(z_f^s, z_f^{p-})\} \end{aligned} \quad (4)$$

where $m^{z_{inv}}, m^{z_f}$ are margin hyperparameters and $\delta(a, b) = \|a - b\|^2$. For simplicity, we have $\mathcal{L}_{\text{Tri}} = \mathcal{L}_{\text{Tri}}^{z_{inv}} + \mathcal{L}_{\text{Tri}}^{z_f}$. Now the overall learning objective of our disentanglement VAE model is:

$$\mathcal{L}_\Omega = \mathcal{L}_{\text{rec}} + \lambda_1 \cdot \mathcal{L}_{\text{KL}} + \lambda_2 \cdot \mathcal{L}_{\text{Tri}}, \quad (5)$$

where $\Omega = \{\phi, \theta\}$; λ_1, λ_2 are weighting hyperparameters.

3.2. Meta-Learning for Adaptive Disentanglement

Overview: Disentangling styles in sketches is more challenging compared to other images like paintings [24]. Besides having sparse visual cues, untrained amateur sketches hold much more variation in style unlike paintings which hold a distinct style-signature being trained under various definite schools of arts. More importantly, the exhibited style even for the same user can vary depending on which object instance depicted. It is thus critical to learn a disentanglement model that is capable of dynamically adapting to any unseen user style as well as object instances. This is achieved through meta-learning.

Task Sampling: In a meta-learning framework [18], a model is trained from various related labelled tasks. To sample a task $T_i \sim p(T)$ here, we first select a random category C_i out of M categories. Out of all n_i sketch-photo pairs in C_i , ‘ r_i ’ randomly chosen pairs are set aside for validation (query) set ($\mathcal{D}_i^{\text{val}}$), while the remaining N_i pairs constitute the training (query) set ($\mathcal{D}_i^{\text{trn}}$). Inner loop update is performed over \mathcal{D}^{trn} with an aim to minimise the loss in the outer loop over \mathcal{D}^{val} . Within every set, hard negatives are chosen from rest $M-1$ categories ensuring completely dissimilar instances. Next, to prepare the VAE described earlier for meta-learning, we introduce two new components.

Meta-enhancing feature encoder: Inspired from [53], the first new component is a set of feature-transformation (FT) layers plugged into the Encoder (Enc_ϕ), with an aim to dynamically minimise the style-variance in sketches. These FT layers are added after the batch-normalisation layers in Enc_ϕ . For an intermediate feature map $\mathcal{F} \in \mathbb{R}^{h' \times w' \times c}$ where h', w' and c are height, width and number of channels respectively, we sample the bias (ω) and scaling (η) terms as : $\omega \sim \mathcal{N}(0, \text{SmoothReLU}(\phi_\omega))$; $\eta \sim$

$\mathcal{N}(1, \text{SmoothReLU}(\phi_\eta))$ where $\phi_\gamma = \{\phi_\omega, \phi_\eta\} \in \mathbb{R}^{1 \times 1 \times c}$ are hyper-parameters that signify the standard deviation of Gaussian distributions for sampling affine transformation parameters. The activation thus changes to: $\hat{F} = \eta \times F + \omega$. As determining hyper-parameters ϕ_γ empirically for every layer across different sketch-styles would be costly, we optimise them via episodic training – commonly adopted in meta-learning. Each training episode consist an inner loop and an outer one. In the inner loop, the model is updated with a training loss:

$$\Omega' \leftarrow \Omega - \alpha \nabla_{\Omega} \underbrace{\mathcal{L}_{trn}(\{\phi, \phi_\gamma\}, \theta; \mathcal{D}^{trn})}_{Enc}, \quad (6)$$

where α is the inner-loop learning rate. Then in the outer loop, the layers are pulled out and a loss (\mathcal{L}_{tst}) is calculated on the validation set using the modified parameters Ω' . As \mathcal{L}_{tst} denotes the efficiency of feature-transformation layer, we update ϕ_γ in the outer-loop (with a learning rate β):

$$\phi_\gamma \leftarrow \phi_\gamma - \beta \nabla_{\Omega'} \sum_{\mathcal{T}_i} \mathcal{L}_{tst}(\Omega'; \mathcal{D}^{val}). \quad (7)$$

Meta-regularising Disentanglement: To adapt the necessary extent of disentanglement, as the second new component, we introduce an episodic regularisation of the disentangled modal-invariant latent representation z_{inv} across tasks [2]. Here the regulariser is denoted as $Reg(\cdot)$ that applies ℓ_1 norm regularization to each of the parameters Ω_{inv} of z_{inv} . In each training episode, a regularisation loss is imposed over the parameters Ω_{inv} of z_{inv} , as

$$\mathcal{L}_{reg} = Reg_\psi(\Omega_{inv}) = \sum_h \psi^{(h)} |\Omega_{inv}^{(h)}|. \quad (8)$$

\mathcal{L}_{reg} is added to the task loss that contributes to the inner loop update of the model i.e $\Omega' \leftarrow \Omega$. In the outer loop, the loss is calculated with updated parameters of Ω'_{inv} which therefore reflects the usefulness of the current regulariser. Consequently, its parameter ψ is updated by \mathcal{L}_{tst} as

$$\psi \leftarrow \psi - \beta \nabla_{\Omega'} \sum_{\mathcal{T}_i} \mathcal{L}_{tst}(\Omega'; \mathcal{D}^{val}). \quad (9)$$

This weighted ℓ_1 loss denotes a learnable weight control mechanism, which adaptively modulates the proportion of semantic knowledge to be retained in Ω_{inv} for efficient disentanglement of the invariant semantic. As the same regulariser is trained across varying tasks in a meta-training paradigm, it is learnt to generalise onto any unseen task characterised by a new style for object category/instance.

Meta-Optimisation: We summarise the overall meta-optimisation objective here from all the learning objectives discussed so far. Following Eq. 5 the model parameters Ω are updated to Ω' in the inner loop with overall meta-training loss as:

$$\begin{aligned} \mathcal{L}_{trn} &= \mathcal{L}_{rec} + \lambda_1 \cdot \mathcal{L}_{KL} + \lambda_2 \cdot \mathcal{L}_{Tri} + \lambda_3 \cdot \mathcal{L}_{Reg}, \\ \Omega' &\leftarrow \Omega - \alpha \nabla_{\Omega} \underbrace{\mathcal{L}_{trn}(\{\phi, \phi_\gamma\}, \theta, \psi; \mathcal{D}^{trn})}_{Enc}. \end{aligned} \quad (10)$$

With updated model parameters, a validation loss is computed over validation set (\mathcal{D}^{val}). Here the meta-learning pipeline is trained alongside regularisation and feature transformation losses to optimise a combined loss. The optimisation objective for the outer loop is thus formulated as:

$$\begin{aligned} \mathcal{L}_{tst} &= \mathcal{L}_{rec} + \lambda_1 \cdot \mathcal{L}_{KL} + \lambda_2 \cdot \mathcal{L}_{Tri} \\ \operatorname{argmin}_{\Omega, \psi, \phi_\gamma} \quad &\mathcal{L}_{tst}(\Omega', \psi, \phi_\gamma; \mathcal{D}^{val}). \end{aligned} \quad (11)$$

As Ω' depends on Ω , ψ and ϕ_γ via inner-loop update (Eq. 10), a higher order gradient needs to be calculated for outer loop optimisation. Notably, the model updates by averaging gradient over meta-batch size of B sampled tasks.

4. Experiments

Datasets: For category-level SBIR, two datasets are used. Following [30, 63], the first dataset used is Sketchy [45] (extended) which contains 75k sketches across 125 categories with about 73k images [30] in total. For the second dataset, sketches are taken from the TU-Berlin Extension [13] which contains 250 object categories with 80 free-hand sketches per category. We further use 204,489 extended natural photo images of the same 250 TU-Berlin categories provided by [62] to construct the photo part of the dataset. For both datasets we split photos from each category as 70 : 10 : 20 for meta-training (N_i), meta-validation (r_i) and retrieval evaluation respectively, with the sketches split into the three sets in the same proportion. Note that there is no overlapping between the three sets, meaning that sketch-styles used in evaluation are not seen during training. For FG-SBIR, two publicly available datasets, QMUL-Chair-V2 and QMUL-Shoe-V2 [60, 44] are used. They contain 2000 (400) and 6730 (2000) sketches (photos) respectively. Out of the photos, we keep 275 (100) for retrieval evaluation and use the rest for training, with 1150 (200) as meta-train and 575 (100) for meta-validation from QMUL ShoeV2 (ChairV2) datasets respectively. As both contain multiple sketches per photo-instance (we choose those photos having at least 3 sketches while training), they suit well to our motivation of modelling the diversity in sketch-styles. The input images (sketch/photo) were resized to 256×256 and 299×299 for SBIR and FG-SBIR respectively.

Implementation Details: We implement our model in PyTorch on a 12GB TitanX GPU. We use InceptionV3 [51] as our encoder network. The decoder architecture consists of a series of stride-2 convolutions with BatchNorm-Relu activation applied to every convolutional layer except in the output which has tanh for activation. The feature extracted from the encoder is projected into three 64 dimensional vectors signifying μ , $\log \sigma^2$ and z_{inv} . In practice, while training we first warm up our basic cross-modal framework (§3.1) for 20 epochs, before inserting $Reg_\psi(\cdot)$ and FT-layers for meta-optimisation (Eq. 10, 11). We use Adam optimiser in both inner and outer loops with learning

Table 1. Comparative results of our model against other methods on FG-SBIR (D → disentanglement baselines).

Methods	Chair-V2		Shoe-V2		
	acc.@1	acc.@10	acc.@1	acc.@10	
SOTA	Triplet-SN [60]	47.65	84.24	28.71	71.56
	Triplet-Attn [49]	53.41	87.56	31.74	75.78
	Triplet-RL [5]	56.54	89.61	34.10	78.82
	CC-Gen [35]	54.21	88.23	33.80	77.86
D	D-TVAE [21]	49.37	81.63	27.62	70.32
	D-DVML [29]	52.78	85.24	32.07	76.23
Others	B-Basic-SN	49.58	85.41	29.45	72.83
	B-SN-Group	50.35	88.28	30.14	75.62
	B-Cross-Modal [50]	52.24	86.58	31.18	73.51
	B-Meta-SN	53.57	87.69	32.74	76.92
Proposed	62.86	91.14	36.47	81.83	

Table 2. Comparative results of our model against other methods on SBIR (D → disentanglement baselines).

Methods	Sketchy (ext)		TU Berlin (ext)		
	mAP	P@200	mAP	P@200	
SOTA	DSH (64 bit) [25]	0.711	0.858	0.521	0.655
	GDH (64 bit) [63]	0.810	0.894	0.690	0.728
D	D-TVAE [21]	0.695	0.839	0.507	0.643
	D-DVML [29]	0.785	0.891	0.648	0.693
Others	B-Basic-SN	0.715	0.861	0.531	0.659
	B-SN-Group	0.738	0.872	0.572	0.661
	B-Cross-Modal [50]	0.763	0.884	0.622	0.688
	B-Meta-SN	0.824	0.897	0.674	0.715
Proposed	0.905	0.927	0.778	0.795	

rates of 0.0005 and 0.0001 respectively. Hyperparameters $\lambda_{1 \rightarrow 3}$ (determined empirically) are set to 0.001, 1.0, 0.7 respectively while λ_1 is increased with linear scheduling to 1.8 for the last 75 of 200 epochs, for better training stability. We use a meta-batch size of 16 and set $\mu^{z_{inv}}$ and μ^{z_f} to 0.5 and 0.3 respectively (further details in supplementary).

Evaluation: Category-level SBIR is evaluated similar to [30] using mean average precision(MAP) and precision at top-rank 200 (P@200). For FG-SBIR we use top-q (acc@q) accuracy. We also design an unconventional metric solely for qualitative comparison of modelling sketch diversity in FG-SBIR. Out of ‘m’ photos with ‘k’ sketches per photo (p_i), we define average retrieval rank $\mathcal{R}_{avg} = \frac{1}{m} \sum_i \mathcal{R}_i$ where \mathcal{R}_i is the rank of retrieving ‘ p_i ’ against sketch ‘ s_i ’. $\forall p_i$, let rank variance $\mathcal{V}_i = \text{variance}(\mathcal{R}_1^i, \mathcal{R}_2^i, \dots, \mathcal{R}_k^i)$ where $\mathcal{R}_j^i |_{j=1}^k$ is the retrieval rank of p_i against its j^{th} sketch-style. Accordingly, average rank variance $\mathcal{V}_{avg} = \frac{1}{m} \sum_i \mathcal{V}_i$. Lower the value of \mathcal{R}_{avg} and \mathcal{V}_{avg} , higher is the score and consistency in retrieval accuracy against varying styles per photo, respectively.

4.1. Competitors

For both category-level and FG-SBIR, we evaluate our method against existing state-of-the-art (SOTA) SBIR methods, and a few relevant latent representation disentanglement baselines adapted to our problem. These include:

(a) SOTA: *Triplet-SN* [60] use Sketch-A-Net as baseline feature extractor trained using triplet loss. *Triplet-Attn-SN*

[49] extended [60] with spatial attention using a higher order HOLEF ranking loss. *CC-Gen* [35] takes a cross-category (CC) domain-generalisation approach, modelling a universal manifold of prototypical visual sketch traits that dynamically embeds sketch and photo, to generalise for unseen categories. *Triplet-RL* [5] leverages triplet-loss based pre-training, followed by RL based fine-tuning for on-the-fly FG-SBIR. We report its results only on completed sketches as early retrieval is not our goal. *DSH* [25] unifies discrete binary code learning with visual sketch/photo feature maps to alleviate geometric distortion between sketches and photos. *GDH* [63] learns a domain-migration network using binary hash codes in a generative adversarial paradigm with cycle-consistency losses, without relying on pixel-level alignment between cross-modal pairs.

(b) Disentanglement methods: None of the existing SOTA SBIR models consider latent representation disentanglement. We therefore choose a number representative disentanglement methods and adapt them for (FG-)SBIR for a fair comparison. For these models, encoded features are used for distance-based retrieval during evaluation as done in our model. *D-TVAE* [21] uses a standard VAE training paradigm with single modality translation. A triplet loss is imposed on the extracted mean feature to bring the sketch and matching photo feature closer while distancing the negative photo-feature. *D-DVML* [29] employs a VAE framework with same-modality translation. It involves disentangling sketch features into invariant and variant components but the disentanglement operation is unregulated unlike ours. Besides VAE losses, the model is also guided by triplet loss objective between the invariant component of the sketch, its matching photo, and its non-matching photo.

(c) Other relevant Baselines: *B-Basic-SN* is a naively built Siamese baseline similar to *Triplet-SN* which replaces its Sketch-a-Net with Inception-V3 as a backbone feature extractor. *B-Cross-Modal* [50] learns a cross-modal latent space in a VAE framework, involving translation amongst multiple (sketch and photo) modalities without disentanglement. We further impose a Triplet loss [60] on the generated latent feature bringing sketch and matching photo features closer while distancing the negative one. *B-Group-SN* is similar to *Triplet-SN* with Inception V3 [44] backbone feature extractor, where we concatenate feature embedding of three sketches against one corresponding category (SBIR) or photo-instance (FG-SBIR) and pass them through a linear layer to match the embedding dimension of the photo. This ensures that the sketch representation holds knowledge on multiple sketch-styles per object to some extent. *B-Meta-SN* simply employs vanilla MAML [14] on a *Triplet-SN* with an Inception-V3 backbone, without the FT layers or any disentangling regulariser. It adapts using inner loop updates across retrieval tasks over categories in SBIR and over instances in FG-SBIR frameworks

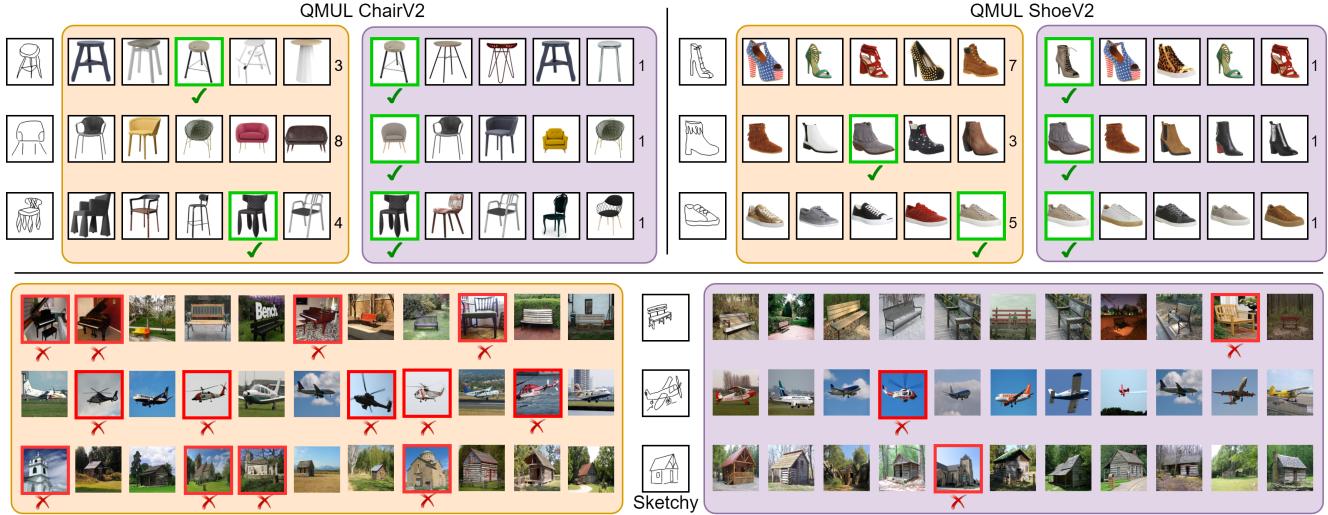


Figure 3. Qualitative retrieval results on QMUL ShoeV2, ChairV2 and Sketchy datasets. B-Basic-SN (orange) vs Ours (magenta).

4.2. Performance Analysis

Comparative results on category-level and FG-SBIR are shown in Table 2 and Table 1 respectively. The following observations can be made: (i) Our method outperforms all other compared methods under both settings and on all four datasets. This clearly illustrates the efficacy of the proposed method thanks to its ability of dynamically adapting to new user styles. (ii) The inferior results of *Triplet-SN* and *Triplet-Attn* are partially due to their apparently weaker backbone feature extractor of Sketch-A-Net. (iii) *Triplet-RL* performs much better owing to its novel reward function designed in reinforcement setup towards sketch completion. *CC-Gen* on the other hand comes close in performance, owing to its learning of universal manifold of visual traits aiding its generalising ability. However both *Triplet-RL* and *CC-Gen* ignore the style diversity issue. Compared to our model, their accuracy is lower by 6.32(2.37)% and 8.65(2.67)% for ChairV2 (ShoeV2) datasets respectively. (iv) For category-level SBIR, although both *GDH* and *DSH* perform well, they are clearly inferior to our method, as they rely on only a singular sketch-image embedding via shared hashing network, without incorporating diverse sketch-styles belonging to the same object. (v) *D-TVAE* performs poorly as it uses the same mean feature used for reconstruction as the modal-invariant component, thus offering sub-optimal disentanglement. In contrast, *D-DVML* fares better owing to better formulated guiding objectives and better modelling of the invariant component of a sketch/photo with higher discriminative knowledge instilled into the model. (vi) Being trained in a *learning-to-learn* setup *B-Meta-SN* performs better than its simpler counterpart *B-Basic-SN* by 3.99 (3.29)%, but lags behind ours by 9.29 (3.73)%, as neither does it disentangle the stylistic variance, nor does it enforce style agnostic encoding of sketches. (vii) *B-Cross-Modal* fares better than *D-*

TVAE as it harnesses greater information owing to learning a latent space aware of cross-modal knowledge. Without learning a disentangled feature space or style-agnostic encoding however, it performs poorer than our model. (viii) *B-SN-Group* performs higher than *Triplet-SN* with a boosted Acc@10 by 4.04(4.06)% as it now holds a stronger understanding of the search space with increased query knowledge. However, without disentanglement and meta-learning for better generalisation, it lags far behind our method.

Diving deeper into the diversity modelling capability of our method, we plot the respective \mathcal{R}_{avg} and \mathcal{V}_{avg} values for some baselines (**B**) on QMUL ChairV2 and ShoeV2 datasets, obtained via our novel metric (§4-Evaluation) in Fig. 4. While the basic siamese net *B-Basic-SN* shows a large variance among retrieval ranks of the same photo using its different sketches, our method has a lower rank variance in addition to a much lower average rank. This proves our method indeed models sketch-style diversity to a considerable extent, thus ensuring higher consistency in retrieval accuracy. Qualitative results are shown in Fig. 3.

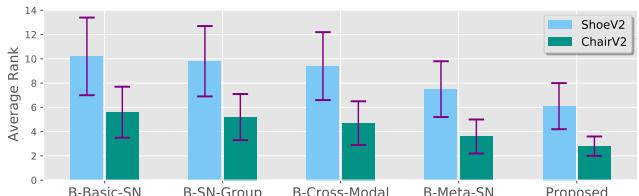


Figure 4. Figure shows proposed method to clearly surpass (lower is better) other baselines in both \mathcal{R}_{avg} (bar height) and \mathcal{V}_{avg} (variance line) in QMUL ShoeV2 and ChairV2 datasets.

4.3. Ablation Study

Is modelling sketching diversity beneficial? For an in depth analysis we perform three experiments: (i) A simple FG-SBIR baseline is trained similar to *Triplet-SN*, on photos and two out of three styles of sketches per photo, keep-

ing the third randomly chosen style per photo for evaluation. Here, the model is guided by triplet loss alone without disentanglement of the learned embedding space for dealing with style diversity. We show the results in Table 3 (**w/o Diversity**) on QMUL-ShoeV2 dataset [44] as it has at least three sketches per photo. Results show that this model performs much poorer than our full model, showing the importance of style diversity modelling through explicit latent representation disentanglement. Some qualitative results with three different styles of the same object (unseen in training) are shown in Fig. 5. (ii) We hypothesise that true efficiency for our sketch representation can be judged in a sketch based sketch retrieval problem. Accordingly we train our model in same meta-learning paradigm in a single modality translation setup (no images) on QuickDraw dataset [16] with 345 categories, where we employ cross-style translation between sketches of same category, and triplet objectives similar to ours, with query sketch (s), its matching sketch (p+) and an unmatched (p-) one. The encoder with its modal-invariant semantic is used for retrieval. Following [57], we use 10k sketches for each category with 8:1:1 split for meta-training, meta-test and evaluation respectively. An observed mAP (P@200) score of **0.748 (0.792)** surpassing [57] by 0.096 (0.103) shows credibility of our method. (iii) Further validating the generality of our sketch-specific design, we perform sketch classification on Quickdraw dataset. We treat our style-agnostic encoder as a feature extractor and fine tune it with a cascaded linear classification layer. A decent score of **75.81%** against Sketch-A-Net’s [61] 68.71%, justifies the generality of our style-agnostic disentangled sketch representation.

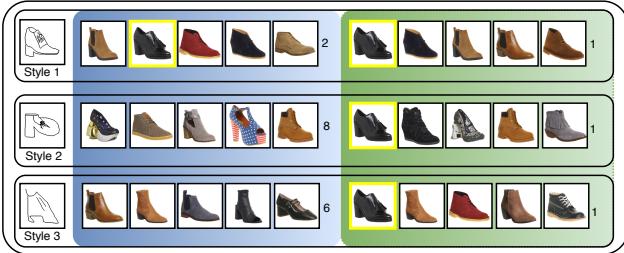


Figure 5. Our (green) method’s efficiency over *w/o-Diversity* (blue) against different styles of the same shoe (highlighted in yellow) is shown. Numbers denote rank of the matching photo.

Importance of Feature-Transformation (FT) Layers: The reason for adding FT layers to the VAE encoder is to impart style-agnostic encoding behavior to Enc_ϕ . To validate this, we train our model without the FT layers (**w/o FT** in Table 3), and observe a dip in performance by 3.19% to 33.28% on ShoeV2. Alternatively, one might empirically set the parameters of FT layers instead of meta-learning it. Accordingly we design a baseline where ϕ_ω and ϕ_η are set to 0.6 and 0.25 respectively for all FT layers (**Fixed-FT**). Inferior results confirm the necessity of meta-learning the hyperparameters ϕ_γ instead of fixing them empirically.

Table 3. Ablative Study

Methods	Shoe-V2		Sketchy (ext)	
	acc.@1	acc.@10	mAP	P@200
w/o Diversity	27.12	69.01	—	—
w/o MFT	33.28	75.34	0.852	0.916
w/o RegD	32.57	73.84	0.837	0.891
Fixed-FT	34.18	79.06	0.878	0.912
Proposed	36.47	81.83	0.905	0.927

Why Regularise Disentanglement?: Our regulariser $Reg_\psi(\cdot)$ modifies the parameters of invariant feature predicting layer z_{inv} to impart the knowledge regarding required extent of disentanglement for adaptively separating the variant and invariant components of sketches. To justify its significance we remove the regularising module and train the model keeping other modules intact (**w/o RegD** in Table 3). We can see that the model’s Acc@1 drops to 32.57% by 3.9% on ShoeV2, which indicates that modelling the diversity in sketches is incomplete without dynamically adapting to the *extent* of disentanglement.

Further Analysis: (i) Following [14] we vary the number of inner loop updates (Eq. 10) during training. Fig. 6 (left) shows a single step update to yield the best performance. Additional updates have a negative impact which might be due to inner loop over-fitting to unnecessary category-specific details of \mathcal{D}^{trn} , thus hampering the generic prior knowledge learned. (ii) Evaluating our framework against varying encoded embedding space dimensions, we find optimal accuracy at $d = 64$, and that performance is stable with higher dimensions (Fig. 6 right). (iii) Having similar evaluation setups, our time-cost (0.18/0.42 ms) for retrieval per query during evaluation on QMUL ChairV2/ShoeV2 lies close to that of *B-Basic-SN* (0.14/0.37 ms).

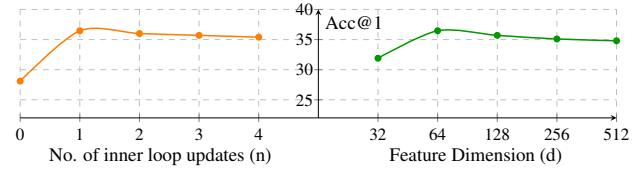


Figure 6. Varying (a) No. of Inner loop updates (optimal at n=1)
(b) feature dimension (optimal at d=64) on QMUL Shoe-V2.

5. Conclusion

In this paper, we addressed a key challenge for sketch-based image retrieval – every person sketches the *same* object *differently*. A novel style-agnostic SBIR model is proposed to explicitly account for the style diversity so that it can generalise onto unseen sketching styles. The model is based on a cross-modal VAE for disentangling a learned latent representation for photo/sketch into a modal-invariant part and a modal-specific part. To make such a disentanglement adaptive to unseen sketch styles, the model is meta-learned with two new components introduced for better generalisation. Extensive experiments show our method to outperform existing alternative approaches significantly.

References

- [1] Antreas Antoniou, Harrison Edwards, and Amos Storkey. How to train your maml. In *ICLR*, 2018. 3
- [2] Yogesh Balaji, Swami Sankaranarayanan, and Rama Chellappa. Metareg: Towards domain generalization using meta-regularization. In *NeurIPS*, 2018. 3, 5
- [3] Ayan Kumar Bhunia, Pinaki Nath Chowdhury, Aneeshan Sain, Yongxin Yang, Tao Xiang, and Yi-Zhe Song. More photos are all you need: Semi-supervised learning for fine-grained sketch based image retrieval. In *CVPR*, 2021. 1
- [4] Ayan Kumar Bhunia, Pinaki Nath Chowdhury, Yongxin Yang, Timothy Hospedales, Tao Xiang, and Yi-Zhe Song. Vectorization and rasterization: Self-supervised learning for sketch and handwriting. In *CVPR*, 2021. 1
- [5] Ayan Kumar Bhunia, Yongxin Yang, Timothy M Hospedales, Tao Xiang, and Yi-Zhe Song. Sketch less for more: On-the-fly fine-grained sketch based image retrieval. In *CVPR*, 2020. 1, 2, 6
- [6] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. In *ICLR*, 2019. 3
- [7] Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *NeurIPS*, 2016. 3
- [8] John Collomosse, Tu Bui, and Hailin Jin. Livesketch: Query perturbations for guided sketch-based visual search. In *CVPR*, 2019. 1, 2
- [9] John Collomosse, Tu Bui, Michael J Wilber, Chen Fang, and Hailin Jin. Sketching with style: Visual search with sketches and aesthetic context. In *ICCV*, 2017. 2, 3
- [10] Sounak Dey, Pau Riba, Anjan Dutta, Josep Llados, and Yi-Zhe Song. Doodle to search: Practical zero-shot sketch-based image retrieval. In *CVPR*, 2019. 1, 2, 11
- [11] Zheng Ding, Yifan Xu, Weijian Xu, Gaurav Parmar, Yang Yang, Max Welling, and Zhuowen Tu. Guided variational autoencoder for disentanglement learning. In *CVPR*, 2020. 3
- [12] Anjan Dutta and Zeynep Akata. Semantically tied paired cycle consistency for zero-shot sketch-based image retrieval. In *CVPR*, 2019. 1, 2
- [13] Mathias Eitz, James Hays, and Marc Alexa. How do humans sketch objects? *ACM TOG*, 2012. 5
- [14] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017. 2, 3, 6, 8
- [15] Chelsea Finn, Aravind Rajeswaran, Sham Kakade, and Sergey Levine. Online meta-learning. In *ICML*, 2019. 3
- [16] David Ha and Douglas Eck. A neural representation of sketch drawings. In *ICLR*, 2017. 1, 8
- [17] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *ICLR*, 2017. 3
- [18] Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural networks: A survey. *arXiv preprint arXiv:2004.05439*, 2020. 3, 4
- [19] Rui Hu and John Collomosse. A performance evaluation of gradient field hog descriptor for sketch based image retrieval. *CVIU*, 2013. 2
- [20] Forrest Huang, John F Canny, and Jeffrey Nichols. Swire: Sketch-based user interface retrieval. In *CHI*, 2019. 2
- [21] Haque Ishfaq, Assaf Hoogi, and Daniel Rubin. Tvae: Triplet-based variational autoencoder using metric learning. 2018. 6
- [22] Zi-Hang Jiang, Qianyi Wu, Keyu Chen, and Juyong Zhang. Disentangled representation learning for 3d face shape. In *CVPR*, 2019. 3
- [23] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR*, 2014. 2, 3
- [24] Dmytro Kotovenko, Artsiom Sanakoyeu, Sabine Lang, and Bjorn Ommer. Content and style disentanglement for artistic style transfer. In *CVPR*, 2019. 2, 4
- [25] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain generalization. In *ICCV*, 2017. 6
- [26] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Learning to generalize: Meta-learning for domain generalization. In *AAAI*, 2017. 3
- [27] Yi Li, Timothy M Hospedales, Yi-Zhe Song, and Shaogang Gong. Fine-grained sketch-based image retrieval by matching deformable part models. In *BMVC*, 2014. 2
- [28] Zhenguo Li, Fengwei Zhou, Fei Chen, and Hang Li. Metasgd: Learning to learn quickly for few-shot learning. *arXiv preprint arXiv:1707.09835*, 2017. 3
- [29] Xudong Lin, Yueqi Duan, Qiyuan Dong, Jiwen Lu, and Jie Zhou. Deep variational metric learning. In *ECCV*, 2018. 6
- [30] Li Liu, Fumin Shen, Yuming Shen, Xianglong Liu, and Ling Shao. Deep sketch hashing: Fast free-hand sketch-based image retrieval. In *CVPR*, 2017. 1, 2, 5, 6
- [31] David G Lowe. Object recognition from local scale-invariant features. In *ICCV*, 1999. 2
- [32] Michael F Mathieu, Junbo Jake Zhao, Junbo Zhao, Aditya Ramesh, Pablo Sprechmann, and Yann LeCun. Disentangling factors of variation in deep representation using adversarial training. In *NeurIPS*, 2016. 3
- [33] Umar Riaz Muhammad, Yongxin Yang, Timothy Hospedales, Tao Xiang, and Yi-Zhe Song. Goal-driven sequential data abstraction. In *ICCV*, 2019. 1
- [34] Xuesong Niu, Zitong Yu, Hu Han, Xiaobai Li, Shiguang Shan, and Guoqing Zhao. Video-based remote physiological measurement via cross-verified feature disentangling. *arXiv preprint arXiv:2007.08213*, 2020. 3
- [35] Kaiyue Pang, Ke Li, Yongxin Yang, Honggang Zhang, Timothy M Hospedales, Tao Xiang, and Yi-Zhe Song. Generalising fine-grained sketch-based image retrieval. In *CVPR*, 2019. 1, 2, 6
- [36] Kaiyue Pang, Yi-Zhe Song, Tony Xiang, and Timothy M Hospedales. Cross-domain generative learning for fine-grained sketch-based image retrieval. In *BMVC*, 2017. 2

- [37] Kaiyue Pang, Yongxin Yang, Timothy M Hospedales, Tao Xiang, and Yi-Zhe Song. Solving mixed-modal jigsaw puzzle for fine-grained sketch-based image retrieval. In *CVPR*, 2020. 2
- [38] Xi Peng, Xiang Yu, Kihyuk Sohn, Dimitris N Metaxas, and Manmohan Chandraker. Reconstruction-based disentanglement for pose-invariant face recognition. In *ICCV*, 2017. 3
- [39] Fabio Pizzati, Pietro Cerri, and Raoul de Charette. Model-based occlusion disentanglement for image-to-image translation. In *ECCV*, 2020. 3
- [40] Yonggang Qi, Yi-Zhe Song, Tao Xiang, Honggang Zhang, Timothy Hospedales, Yi Li, and Jun Guo. Making better use of edges via perceptual grouping. In *CVPR*, 2015. 2
- [41] Andrei A Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. In *ICLR*, 2019. 3
- [42] Jose M Saavedra. Sketch based image retrieval using a soft computation of the histogram of edge local orientations (s-helo). In *ICIP*, 2014. 2
- [43] Jose M Saavedra, Juan Manuel Barrios, and S Orand. Sketch based image retrieval using learned keyshapes (lks). In *BMVC*, 2015. 2
- [44] Aneeshan Sain, Ayan Kumar Bhunia, Yongxin Yang, Tao Xiang, and Yi-Zhe Song. Cross-modal hierarchical modelling for fine-grained sketch based image retrieval. In *BMVC*, 2020. 1, 2, 3, 5, 6, 8
- [45] Patsorn Sangkloy, Nathan Burnell, Cusuh Ham, and James Hays. The sketchy database: learning to retrieve badly drawn bunnies. *ACM TOG*, 2016. 2, 5
- [46] J. Snell, K. Swersky, and R. S. Zemel. Prototypical networks for few shot learning. In *NeurIPS*, 2017. 3
- [47] Jifei Song, Kaiyue Pang, Yi-Zhe Song, Tao Xiang, and Timothy M Hospedales. Learning to sketch with shortcut cycle consistency. In *CVPR*, 2018. 1, 2
- [48] Jifei Song, Yi-Zhe Song, Tony Xiang, and Timothy M Hospedales. Fine-grained image retrieval: the text/sketch input dilemma. In *BMVC*, 2017. 2
- [49] Jifei Song, Qian Yu, Yi-Zhe Song, Tao Xiang, and Timothy M Hospedales. Deep spatial-semantic attention for fine-grained sketch-based image retrieval. In *ICCV*, 2017. 1, 2, 6
- [50] Adrian Spurr, Jie Song, Seonwook Park, and Otmar Hilliges. Cross-modal deep variational hand pose estimation. In *CVPR*, 2018. 4, 6
- [51] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016. 5
- [52] Giorgos Tolias and Ondrej Chum. Asymmetric feature maps with application to sketch based retrieval. In *CVPR*, 2017. 2
- [53] Hung-Yu Tseng, Hsin-Ying Lee, Jia-Bin Huang, and Ming-Hsuan Yang. Cross-domain few-shot classification via learned feature-wise transformation. In *ICLR*, 2020. 3, 4
- [54] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, and Daan Wierstra. Matching networks for one shot learning. In *NeurIPS*, 2016. 3
- [55] Yitong Wang, Dihong Gong, Zheng Zhou, Xing Ji, Hao Wang, Zhifeng Li, Wei Liu, and Tong Zhang. Orthogonal deep features decomposition for age-invariant face recognition. In *ECCV*, 2018. 3
- [56] Lingyu Wei, Liwen Hu, Vladimir Kim, Ersin Yumer, and Hao Li. Real-time hair rendering using sequential adversarial networks. In *ECCV*, 2018. 3
- [57] Peng Xu, Yongye Huang, Tongtong Yuan, Kaiyue Pang, Yi-Zhe Song, Tao Xiang, Timothy M Hospedales, Zhanyu Ma, and Jun Guo. Sketchmate: Deep hashing for million-scale human sketch retrieval. In *CVPR*, 2018. 8
- [58] Linlin Yang and Angela Yao. Disentangling latent hands for image synthesis and pose estimation. In *CVPR*, 2019. 3
- [59] Sasi Kiran Yelamarthi, Shiva Krishna Reddy, Ashish Mishra, and Anurag Mittal. A zero-shot framework for sketch based image retrieval. In *ECCV*, 2018. 2
- [60] Qian Yu, Feng Liu, Yi-Zhe Song, Tao Xiang, Timothy M Hospedales, and Chen-Change Loy. Sketch me that shoe. In *CVPR*, 2016. 2, 4, 5, 6
- [61] 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. *IJCV*, 2017. 8
- [62] Hua Zhang, Si Liu, Changqing Zhang, Wenqi Ren, Rui Wang, and Xiaochun Cao. Sketchnet: Sketch classification with web images. In *CVPR*, 2016. 5
- [63] Jingyi Zhang, Fumin Shen, Li Liu, Fan Zhu, Mengyang Yu, Ling Shao, Heng Tao Shen, and Luc Van Gool. Generative domain-migration hashing for sketch-to-image retrieval. In *ECCV*, 2018. 2, 5, 6
- [64] Yang Zou, Xiaodong Yang, Zhiding Yu, BVK Kumar, and Jan Kautz. Joint disentangling and adaptation for cross-domain person re-identification. In *ECCV*, 2020. 3

Supplementary material for StyleMeUp: Towards Style-Agnostic Sketch-Based Image Retrieval

Aneeshan Sain^{1,2} Ayan Kumar Bhunia¹ Yongxin Yang^{1,2}
Tao Xiang^{1,2} Yi-Zhe Song^{1,2}

¹ SketchX, CVSSP, University of Surrey, United Kingdom

² iFlyTek-Surrey Joint Research Centre on Artificial Intelligence.

Additional explanations

Clarity on bridging domain gap:

From the viewpoint of bridging the domain gap, a gradient reversal layer is employed in Dey *et al.* [10], that is used to create a domain-agnostic embedding, which however does not differentiate if it comes from a sketch or a photo. Our motivation is different – in addition to tackling the sketch-photo domain gap, we further focus on narrowing the domain gaps that exist amongst different sketching styles (i.e., learning a style-agnostic embedding). In particular, the *feature transformation layer* helps bridge this style gap by simulating varying distributions in the intermediate layers of the encoder, and thus condition the encoder to *generalise onto unseen sketching styles*. The meta-learning paradigm further ensures that this notion of style variance is minimised over episodic training, finally resulting in a style-agnostic embedding.

Additional experimental comparison:

The results of DSH and GDH on Sketchy and TU-Berlin have been taken directly from their respective papers. For further transparency we re-run these baselines using Inception-V3 as backbone. Table 4 shows these results to be in line with our conclusions for Sketchy and TUBerlin datasets respectively –

Table 4. Quantitative analysis using Inception-V3 backbone

Method	Sketchy		TUBerlin	
	mAP	P@200	mAP	P@200
DSH	0.725	0.867	0.537	0.660
GDH	0.821	0.896	0.696	0.741
Ours	0.905	0.927	0.778	0.795

More on training details:

The hyperparameters $\lambda_{1 \rightarrow 3}$ have been determined empirically. The impact of \mathcal{L}_{KL} is suppressed ($\lambda_1=0.001$) during initial stages of training, and increased with linear scheduling later for better training stability. We further observed that λ_2 works best if kept constant throughout. Changing λ_3 had generally produced comparatively lower results. Margin hyperparameters for triplet losses $\mu^{z_{inv}}$ and μ^{z_f} were set empirically as well. Please note that unlike few-shot adaption in MAML, there is no adaptation step here during inference. Instead, meta-learning is employed only during training to learn a style-agnostic feature encoder for better generalisation.

More on Fusing modal invariant and modal specific features:

Combining these two components helps the model in keeping important details that might have been removed during disentanglement, for image (sketch/photo) reconstruction. Furthermore, as we intend to learn *how to disentangle* modal-invariant feature from modal-specific one, combining them to obtain a proper reconstruction re-verifies that the disentanglement itself has been learned properly. However, experimental results suggested that element-wise addition performs better than concatenating the two components together. This is probably because the former establishes a clearer boundary between the disentangled components than concatenation.