

Style-Content Disentanglement in Language-Image Pretraining Representations for Zero-Shot Sketch-to-Image Synthesis

Jan Zuiderveld

jan.zuiderveld@student.uva.nl

University of Amsterdam

Royal Conservatoire The Hague

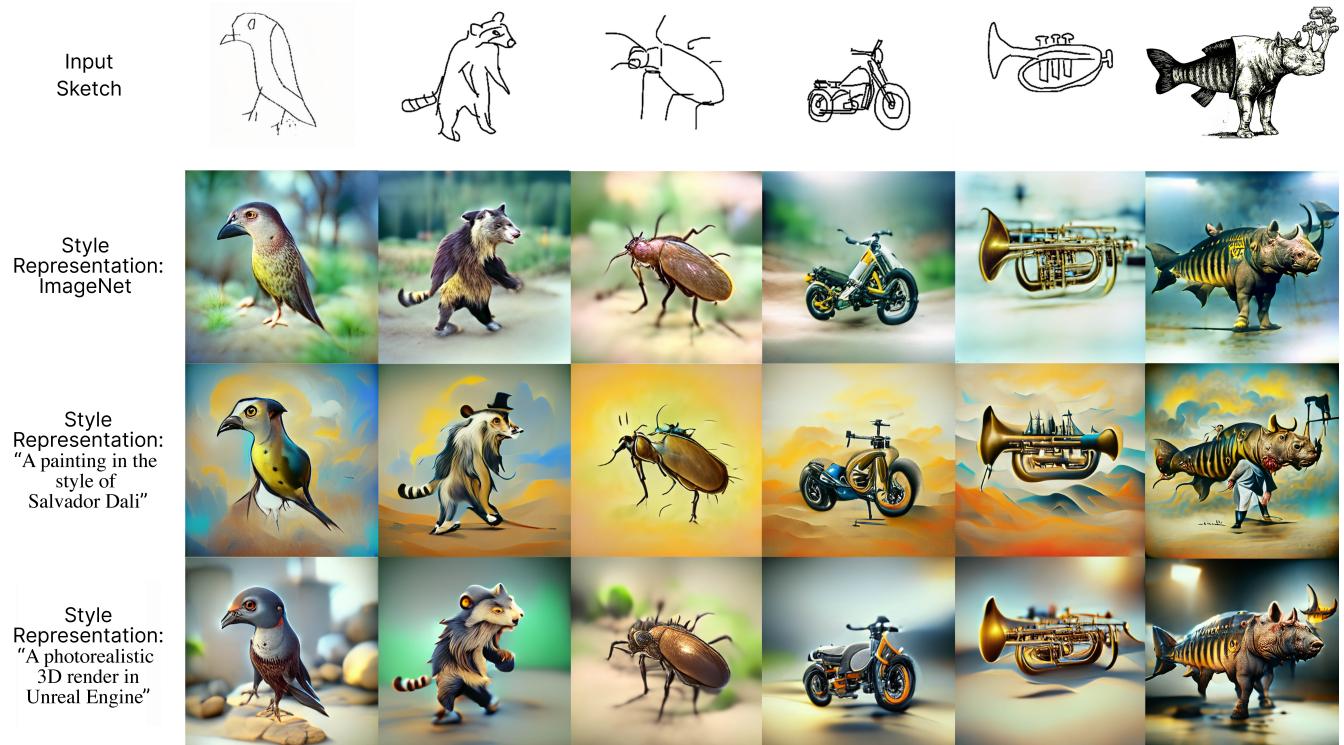


Figure 1: Several results of our proposed framework, with given input sketches and style representations.

ABSTRACT

In this work, we propose and validate a framework to leverage language-image pretraining representations for training-free zero-shot sketch-to-image synthesis.

We show that disentangled content and style representations can be utilized to guide image generators to employ them as sketch-to-image generators without (re-)training any parameters. Our approach for disentangling style and content entails a simple method

consisting of basic arithmetic assuming compositionality of information in representations of input sketches. Our results demonstrate that this approach is competitive with several state-of-the-art object-level domain-specific sketch-to-image models, while only depending on pretrained off-the-shelf models.

We limit our reported experiments to sketch-to-image synthesis, but preliminary results show that our method is applicable to more image-to-image translation domains.

CCS CONCEPTS

- Computing methodologies → Image representations; Image-based rendering.

KEYWORDS

style-content disentanglement, sketch-to-image synthesis, language-image pretraining

117 **ACM Reference Format:**

118 Jan Zuiderveld. 2018. Style-Content Disentanglement in Language-Image
 119 Pretraining Representations for Zero-Shot Sketch-to-Image Synthesis. In
 120 *Proceedings of Association for Computing Machinery's Special Interest Group*
 121 *on Computer Graphics and Interactive Techniques (ACM SIGGRAPH '22)*.
 122 ACM, New York, NY, USA, 4 pages. <https://doi.org/XXXXXXX.XXXXXXX>

123

124 **1 INTRODUCTION**

126 The power and promise of deep generative models such as GANs
 127 lie in their ability to synthesize endless realistic, diverse, and novel
 128 content with minimal user effort. The potential utility of these mod-
 129 els continues to grow thanks to the increased quality and resolution
 130 of large-scale generative models in recent years. Conditional im-
 131 age synthesis allows users to use inputs to control the output of
 132 image synthesis methods. Over the years, a variety of input modal-
 133 ities have been studied, mostly based on conditional GANs. Several
 134 methods exist that allow generative models to be guided by condi-
 135 tioning variables, e.g. text, sketches or segmentation maps. Sketch
 136 conditioning finds a practical middle ground within the character-
 137 istics of these methods. It allows significantly more control over
 138 output structure than text conditioning, and does not require labels
 139 as segmentation conditioning does.

140 However, current sketch conditioned models are trained on a
 141 relatively narrow set of visual concepts. This restricts the space
 142 of possible inputs that result in desired output images, impeding
 143 the general usability of these models. An important cause of this is
 144 the fact that paired image datasets are labor intensive and costly to
 145 create. Approaches that circumvent this requirement exist, these
 146 include CycleGAN based architectures (1909.08313, 2104.05703.pdf)
 147 and the usage of image-to-sketch models for synthesizing paired
 148 data (2012.09290.pdf and gaugan2). But as of yet, no work has
 149 been published applying sketch-to-image on large, highly diverse
 150 datasets.

151 Recently, the idea of learning visual concepts from supervi-
 152 sion contained in natural language, Language-Image Pre-training,
 153 has gained a lot of attention. These models learn unified visual- and
 154 natural language representations from the supervision contained
 155 in the vast amount of text and associated images on the internet,
 156 resulting in very rich visual representations.

157 In this paper, we propose a framework to leverage classifier
 158 representations for sketch-to-image generation. A simple method
 159 consisting of basic arithmetic assuming compositionality of infor-
 160 mation in the representations is used to disentangle spatial and
 161 conceptual information from low-level, stylistic information in rep-
 162 resentations of input sketches. These disentangled representations
 163 can then be utilized to guide image generators to employ them as
 164 sketch-to-image generators without (re-)training any parameters.

165 The results in this paper and the supplementary material demon-
 166 strate that our approach, when used with sufficiently flexible image
 167 generators and informative representations, several state-of-the-art
 168 object-level domain-specific sketch-to-image models, while only
 169 depending on pretrained off-the-shelf models and being evaluated
 170 in a zero-shot fashion. Figure ?? shows several examples using our
 171 approach. The effectiveness of our method indicates that CLIP rep-
 172 resentations exhibit strong compositionality, something that has
 173 not been harnessed for disentanglement in earlier research.

174

175 We limit our reported experiments to sketch-to-image synthe-
 176 sis, but preliminary results show that our method is applicable to
 177 more image-to-image translation domains with diverse content.
 178 See Appendix ?? for a few experiments including photo-to-sketch
 179 and ...

180

181

182 **2 RELATED WORK**183 **2.1 Sketch to Image Synthesis**

184 The goal of sketch-based image synthesis is to output a target image
 185 from a given sketch. Early works [8, 18, 9] regard freehand sketches
 186 as queries or constraints to retrieve each composition and stitch
 187 them into a picture, requiring availability of close image matches
 188 in the queried database for reasonable results. In recent years, an
 189 increasing number of works adopt GAN-based models [21] to learn
 190 sketch-to-image synthesis directly.

191

192 [80, 39, 7] train a Conditional GAN (Mirza and Osindero 2014)
 193 with photos and corresponding edge maps to make up for the lack
 194 of real sketch data, but only report results trained on single class
 195 datasets. However, sketches differ from edge maps in several ways,
 196 resulting in inferior results when applying these models to sketches.
 197 The authors of SketchyGAN [10] gathered the largest currently
 198 available paired dataset of sketches and pictures (75k sketches
 199 across 125 categories), and trained a conditional GAN on these.
 200 ContextualGAN [48] turns the image generation problem into an
 201 image completion problem: the network learns the joint distribution
 202 of sketch and image pairs and acquires the result by iteratively
 203 traversing the manifold. PoE-GAN is a multimodal conditional GAN
 204 that is trained on images with paired text, sketch and segmentation
 205 maps, achieving state-of-the-art sketch-to-image synthesis.

206

207 All of these proposed architectures are trained on datasets with
 208 a limited amount of classes. To the best of our knowledge, there
 209 is no published research reporting training any kind of image-to-
 210 image translation architectures on large-scale, diverse datasets for
 211 category-free, zero-shot image-to-image generation. Our work is
 212 the first to try zero-shot image-to-image generation .

213

214 **2.2 Language-Image Pre-training**

215 Multiple recent works learn cross-modal vision and language rep-
 216 resentations [12, 47, 52, 35, 30, 51, 29, 7, 32] for a variety of tasks,
 217 such as language-based image retrieval, image captioning, and vi-
 218 sual question answering. Following the success of BERT [13] in
 219 various language tasks, recent vision and language methods typ-
 220 ically use Transformers [55] to learn the joint representations. A
 221 recent model, based on Contrastive Language-Image Pre-training
 222 (CLIP) [42], consisting of a transformer- language encoder f_{text} and
 223 image encoder f_{image} is trained to learn a multi-modal representa-
 224 tion space. which can be used to estimate the semantic simila-
 225 rity between a given text \mathcal{T} and an image \mathcal{I} by evaluating their cosine
 226 simialrity

227

228 CLIP was trained on 400 million text-image pairs, collected from
 229 a variety of publicly available sources on the Internet. The represen-
 230 tations learned by CLIP have been shown to be extremely powerful,
 231 enabling state-of-the-art zero-shot image classification on a variety
 232 of datasets.

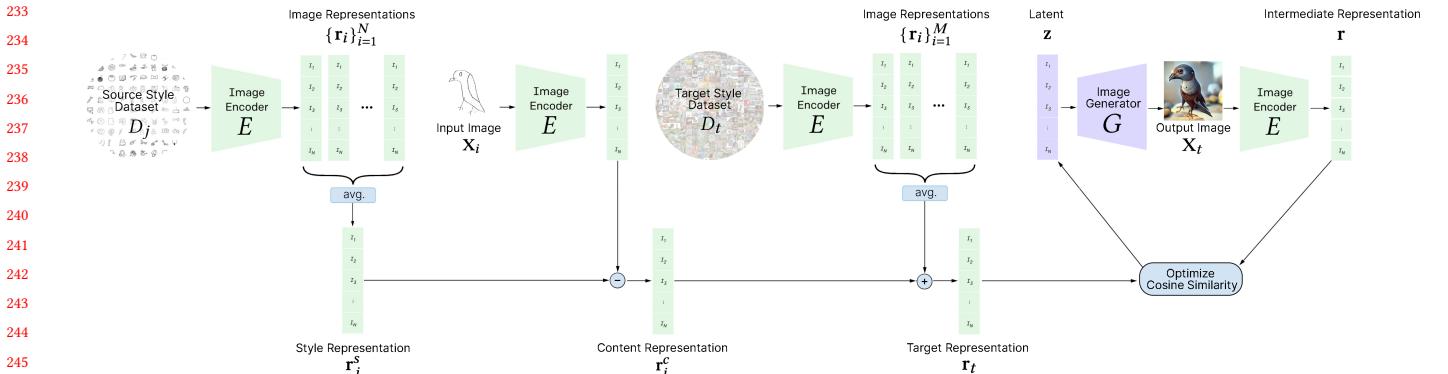


Figure 2: Overview of our proposed framework.

2.3 Classifier Guidance Based Image Synthesis

GANs for conditional image synthesis [39, 5] make heavy use of class labels. This often takes the form of class-conditional normalization statistics [16, 11] as well as discriminator heads that are explicitly designed to behave like classifiers [40], indicating that class information is crucial to the success of these models. The general idea of optimizing latent representations or parameters of an image generator using a separately trained classifier has been widely used as a powerful framework for generating, editing and recovering images; see, e.g., [1, 14, 16, 39, 2105.05233.pdf], to name only a few. Santurkar et al. X show that adversarial robustness of the classifier is a crucial aspect for optimal execution of such tasks.

Several concurrent projects use CLIP as a classifier to guide text-to-image generation through optimization. Deep Daze [38] optimizes the weights of an implicit neural representation network, while [37, 41, 10, 16] optimize the latent space of BigGAN [5], StyleGAN [25] or VQGAN [15]. GLIDE et al. propose a noised CLIP guided diffusion architecture.

3 METHODOLOGY

Much current research in unsupervised image-to-image translation relies on learning content and style representations. In these approaches, the content c_i of input images X_i and the style of target images s_k together condition a generator G that synthesizes translated images $X_t = G(c_i, s_t)$.

Instead of learning content and style representations from scratch, we propose a simple method for disentangling style and content in representations of a pretrained image encoder E under the assumptions of in-distribution input and compositionality of representations. We use obtained representations to guide G . We apply this framework to do sketch-to-image synthesis. See Figure ?? for a general overview.

3.1 Problem Formulation

Let us model the formation of an image X as a function of style $s \in \mathcal{S}$ and content $c \in \mathcal{C}$, $X = G(s, c)$. Let \mathcal{P} denote the set of possible parts and \mathcal{X} the input space. For each $X \in \mathcal{X}$, we assume the existence of a function D mapping X to $\mathcal{P}' \subseteq \mathcal{P}$, the set of its parts. These parts can be any feature of the input, e.g. the presence

and location of an object, features of this object such as its size and color, or global style of an image. Style s describes all parts $p^s \in D(X)$ that are invariant within the stylistic domain of X , while content c describes all residual parts $p^c \in D(X)$ that are not described by style s .

Given a dataset D_j of N images X_i , all sampled from stylistic domain s_j with normally distributed content c_i ,

$$D_j = \{X_i\}_{i=1}^N, \quad X_i = G(s_j, c_i), \quad c_i \sim \mathcal{N}(0, \sigma^2 I)$$

We use image encoder $E: \mathcal{X} \rightarrow \mathcal{R}$ to obtain representations $r_i = E(X_i)$, where \mathcal{R} denotes the representation space. We assume r_i to be compositional, it can be expressed as a weighted sum of simpler parts. Let $h: \mathcal{P} \rightarrow \mathcal{R}$ denote a function that maps parts to representations. Formally, we define a function $f(\mathcal{X}) \in \mathcal{R}$ as compositional if it can be expressed as a weighted sum of the elements of $\{h(p) \mid p \in D(\mathcal{X})\}$ (Brendel Bethge, 2019b).

We decompose r_i as a sum of content representation r_i^c and style representation r_i^s , $r_i = r_i^c + r_i^s$. We assume X_i to be in-distribution for $E(\mathcal{X})$, thus r_i^c follows the distribution of c_i . We can then obtain r_i^s by taking the arithmetic mean over all r_i ,

$$r_j^s = \frac{1}{N} \sum_{i=1}^N r_i = \frac{1}{N} \sum_{i=1}^N r_i^c + \frac{1}{N} \sum_{i=1}^N r_i^s = 0 + \frac{N}{N} r_j^s.$$

Finally, to obtain r_i^c we subtract r_j^s from r_i ,

$$r_i^c = r_i - r_j^s = r_i^c + r_j^s - r_j^s.$$

Using these operations we obtain r_i^c describing the content of input images X_i from stylistic domain s_j , and r_t^s describing target stylistic domain s_t of dataset D_t . We sum these to create target representations for guiding image synthesis: $r_t = r_i^c + r_t^s$.

Specifically, given image generator $G(z)$, and image encoder E , we solve the following optimization problem to synthesize X_t , a translation of X_i in style s_t :

$$\arg \min_{z \in \mathcal{Z}} \frac{\langle E(G(z)), r_t \rangle}{\|E(G(z))\| \cdot \|r_t\|},$$

where $\langle \cdot, \cdot \rangle$ computes the cosine similarity between its arguments.

349 3.2 Improving Adversarial Robustness

350 Most off-the-shelf classifiers are not trained for exhibiting adversarial robustness. Since we are directly optimizing images using
 351 such classifiers, this process is prone to induce non-semantic perturbations in X_t that increase classification scores [cite fuseddream].
 352 Inspired by self-supervised learning methods [cite], we use augmentation pipeline a when feeding input X_t to image encoder e :

$$353 \quad a(X_t) = \mathbb{E}_{X'_t \sim \pi(\cdot | X_t)} a(X'_t)$$

354 where X'_t is a random perturbation of the input image X_t drawn
 355 from distribution $\pi(\cdot | X_t)$ of candidate data augmentations, including
 356 random colorization, translation, and cutout. This generates
 357 new samples X'_t that average out adversarial gradients while pre-
 358 serving most of the content and style information in X_t .

359 4 EXPERIMENTAL SETUP

360 4.1 Image Encoder E

361 To evaluate the adopted framework we report experiments with
 362 CLIP [cite] as the image encoder E . CLIP is well-suited for our frame-
 363 work as it's large-scale dataset it has been trained on ensures that
 364 a wide range of styles s are in-distribution. And, more importantly,
 365 the excellent zero-shot performance of CLIP indicates its represen-
 366 tations exhibit a significant amount of compositionality, making
 367 it well-suited for our style-content disentanglement method. As
 368 demonstrated by X et al. [1912.12179.pdf], zero-shot classification
 369 performance strongly correlates with compositionality in learned
 370 representations. This is very intuitive: we expect compositionality

371 to be an advantage in zero-shot learning: if a model has a good
 372 understanding of how parts map to representations, it can learn to
 373 combine known concepts to describe new classes.

374 Since CLIP also includes a natural language encoder which maps
 375 text to the same embedding space as its image encoder, we also
 376 experiment with target styles s_t obtained through text encodings.
 377 See Figure ??.

378 4.2 Image Generator G

379 The space of images our approach is limited by the representation
 380 power of the GAN we use. This makes the method difficult to
 381 generate out-of-distribution images and prone to inherit the data
 382 biases e.g., center, spatial and color bias [2, 16], from the original
 383 training set of the GAN.

384 Since we aim to achieve zero-shot sketch-to-image synthesis,
 385 the image generator needs to be flexible. However, very flexible im-
 386 age generators such as differentiable image parameterizations, e.g.
 387 implicit neural representations [cite] or direct pixel optimization,
 388 facilitate adversarial attacks.

389 In section ?? we report results of using VQGAN [cite] in the
 390 adopted framework. this architecture exhibits spatially flexible im-
 391 age synthesis due to its latent embedding structure. We initialize
 392 z by using VQGAN's image encoder to encode input sketch X_i .
 393 Additionally, we report results of using GLIDE (guided by noised
 394 CLIP) [cite] in section ??.

395 REFERENCES