# Unsupervised Scene Sketch to Photo Synthesis

Jiayun Wang[1]          Sangryul Jeon[1]          Stella X. Yu[1]
Xi Zhang[2]          Himanshu Arora[2]          Yu Lou[2]

[1] UC Berkeley / ICSI                    [2] Amazon
{peterwg,srjeon,stellayu}@berkeley.edu          {xizhn,arorah,ylou}@amazon.com

**Abstract.** Sketches make an intuitive and powerful visual expression as they are fast executed freehand drawings. We present a method for synthesizing realistic photos from scene sketches. Without the need for sketch and photo pairs, our framework directly learns from readily available large-scale photo datasets in an unsupervised manner. To this end, we introduce a standardization module that provides *pseudo* sketch-photo pairs during training by converting photos and sketches to a standardized domain, i.e. the edge map. The reduced domain gap between sketch and photo also allows us to disentangle them into two components: holistic scene structures and low-level visual styles such as color and texture. Taking this advantage, we synthesize a photo-realistic image by combining the structure of a sketch and the visual style of a reference photo. Extensive experimental results on perceptual similarity metrics and human perceptual studies show the proposed method could generate realistic photos with high fidelity from scene sketches and outperform state-of-the-art photo synthesis baselines. We also demonstrate that our framework facilitates a controllable manipulation of photo synthesis by editing strokes of corresponding sketches, delivering more fine-grained details than previous approaches that rely on region-level editing.

**Keywords:** sketch, scene sketch, photo synthesis, unsupervised learning

## 1 Introduction

Sketching is an intuitive way to represent visual signals. With a few sparse strokes, humans could understand and envision a photo from a sketch. Additionally, unlike photos which are rich in color and texture, sketches are easily editable as strokes are easy to modify. We aim to synthesize photos that preserve the structure of scene sketches while delivering the low-level visual style of reference photos.

Unlike previous works [15, 24, 32] that synthesize photos from categorical object-level sketches, our goal in which scene-level sketches are used as input poses additional challenges due to **1) Lack of data.** There is no training data available for our task due to the complexity of scene sketches. Not only the insufficient amount of scene sketches, but the lack of paired scene sketch-image datasets make supervised learning from one modality to another intractable. **2) Complexity of scene sketches.** A scene sketch usually contains many
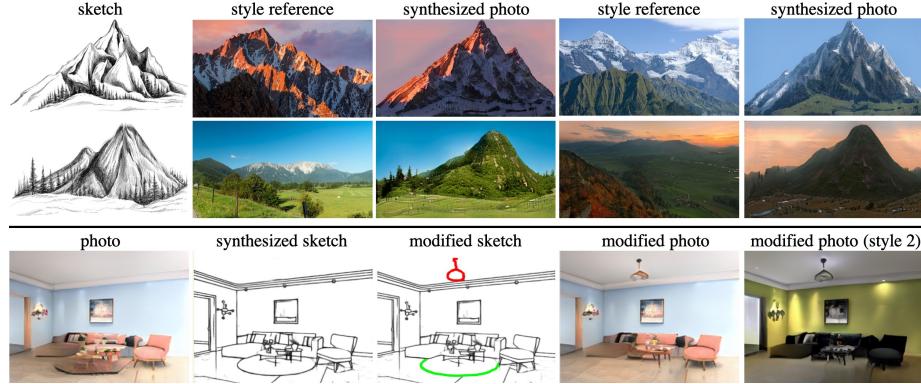
Fig. 1: *Upper:* Given a sketch and a style reference photo, our method is capable of transferring low-level visual styles of the reference while preserving the content structure of the sketch. We show synthesis results with different references. *Lower:* Given an arbitrary photo, users could easily and interactively edit it by adding or removing strokes on the synthesized sketch.

objects of diverse semantic categories with complicated spatial organization and occlusions. Isolating objects, synthesizing object photos and combining them together [7] do not work well and are hard to generalize. For one, detecting objects from sketches is hard due to the sparse structure. For another, one may encounter objects that do not belong to seen categories, and the composition could also make the synthesized photo unrealistic.

We propose to alleviate these issues via **1)** a standardization module, and **2)** disentangled representation learning.

For the lack of data, we propose a standardization module, where input images are converted to a standardized domain, edge maps. Edge maps can be considered as *synthetic sketches* due to the high similarity to real sketches. With the standardization, readily-available large-scale photo datasets could be used for training by converting them to edge maps. Additionally, during inference, sketches of various individual styles are also standardized such that the gap between training and inference is narrowed.

For the complexity of scene sketches, we learn disentangled holistic content and low-level style representations from photos and sketches by encouraging only content representations of photo-sketch pairs to be similar. As a definition, content representations encode holistic semantic and geometric structures of a sketch or photo. Style representations encode the low-level visual information such as color and texture. A sketch could depict similar contents as a photo, but contain no color or texture information. By factorizing out colors and textures, the model could directly learn from large-scale photos for scene structures and transfer the knowledge to sketches. Additionally, combining the content representation of a sketch and a style representation of a reference photo could decode a realistic photo. The decoded photo should depict similar contents as the sketch and shares a similar style with the reference photo. This is the underlying mechanics of the
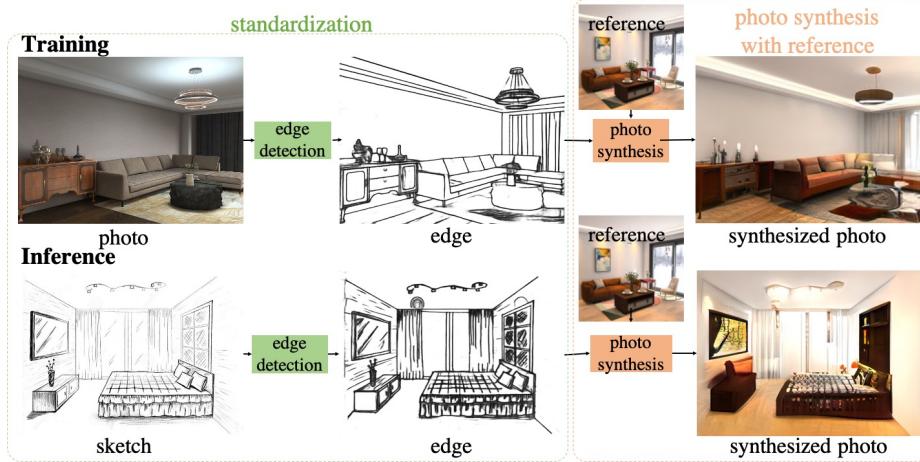
Fig. 2: Our method consists of two components, standardization and photo synthesis. **Left:** The standardization module converts photos or sketches into a standardized domain, edge maps, to reduce the domain gap between training and inference. **Right:** From the standardized edge map, the photo synthesis module generates a photo with a similar style as the given reference image.

proposed reference-guided scene sketch to photo synthesis approach. Note that the disentangled representations have been studied previously for photos [28, 34] and we extend the concept to sketches.

As exemplified in Fig.1, not only photo synthesis from scene sketch, our model can promote also controllable photo editing by allowing users to directly modify strokes of a corresponding sketch. The process is easy and fast as strokes are easy and flexible to modify, compared with photo editing from segmentation maps proposed by previous works [15, 22, 26, 28]. Specifically, the standardization module first converts a photo to a sketch. Users could modify strokes of the sketch and synthesize a newly edited photo with our model. Additionally, the style of the photo could also be modified with another reference photo as guidance.

We summarize our contribution as follows: **1)** We propose an unsupervised scene sketch to photo synthesis framework. We introduce a standardization module that converts arbitrary photos to standardized edge maps, enabling a vast amount of real photos to be utilized during training. **2)** Our framework facilitates controllable manipulation of photo synthesis through editing scene sketches with more plausibility and simplicity than previous approaches. **3)** Technically, we propose novel designs for scene sketch to photo synthesis, including shared content representations to enable knowledge transfer from photos to sketches and model fine-tuning with sketch-reference-photo triplets for improved performance.

## 2   Related Work

**Conditional Generative Models**. Previous approaches generated realistic images by conditioning generative adversarial networks [9] on a given input

from users. More recent methods extended it to multi-domain and multi-modal setting [4, 13, 23], facilitating numerous downstream applications including image inpainting [14, 29], photo colorization [20, 40], texture and geometry synthesis [10, 42]. However, naively adopting this framework to our problem is challenging due to the absence of paired data where sketches and photos aligned. We address this by projecting arbitrary sketches and photo into the intermediate representation and generating pseudo paired data to learn in an unsupervised setting.

**Disentanglement of Content and Style Representations.** The disentanglement has been studied [31, 44] prior to the surge of deep learning models, where they show low-level style like texture can be modeled as statistics of an image. Deep generative models [16, 21, 28, 34] also achieved success in photo style transfer by the disentanglement. We extend the disentanglement idea to sketches and show its application in photo synthesis.

**Sketch to Photo Synthesis**. Following a seminal work, SketchGAN [3], several efforts has been made on synthesizing photos [8, 24, 37] or reconstructing 3D shapes [5, 35, 36] from sketches. They however mainly focused on categorical single-object sketches without substantial background clutters, and thus have difficulties when encountered with complicated scene-level sketches.

Scene sketch to photo synthesis is limited by lack of the data. SketchyScene [45] is the only scene dataset with object segmentation and corresponding cartoon images. However, their sketch is manually composited from multiple object sketches with reference to a cartoon image. The composite sketch has a large domain gap to real scene sketches with reference to a real scene. Their composition idea greatly impacts how researchers solve the photo synthesis. [7] detect objects of composite sketches and generate individual photos as well as a background image and combine them together. Holistic scene structures are ignored and the photo composition leads to artifacts and unrealism. We learn holistic scene structures from massive photo datasets and transfer the knowledge to sketches.

**Deep Image Editing**. By the favor of powerful generative models [17], previous works edited photos by modifying the extracted latent vector. Typically they sampled the desired latent vector from a fixed distribution according to a user's semantic control [43], or let a user spatially annotate the region-based semantic layout [27, 28]. DeepFaceDrawing [2] enables user to sketch progressive for face image synthesis. Our work differs in that we allow users to directly edit strokes of a complicated scene sketch, thus enabling much more fine-grained editing.

## 3   Methods

As illustrated in Fig.2, our framework mainly consists of two components: domain standardization and reference-based photo synthesis. For standardization (details in Section 3.1), input photos and sketches are converted to standardized edge maps, which bypass the lack of data issue. The second part is reference-guided photo synthesis (details in Section 3.2), where synthesized photos are generated based on input sketches and style reference photos.

Fig. 3: The standardization module converts photos and sketches to a standardized domain, edge maps. After the standardization, edges of photos and sketches share higher similarity, which makes the domain gap between training and evaluation narrower. Within the test set, edges of sketches with different individual styles also a share higher similarity, making the intra-sketch-set discrepancy smaller.

### 3.1   Domain Standardization

Due to the lack of paired sketch-photo datasets, it is intractable for supervised models to synthesize photos from sketches. We adopt a similar idea as [35], where they converted inputs to a standardized domain, and showed learning from such domain has better performance compared to directly using unprocessed inputs.

As shown in Fig.2**L**, the standardization can be considered as data prepossessing and is different for training and inference. During training, we collect a large scale photo dataset of a specific category, e.g., indoor scenes. Each photo is converted to a standardized edge map for later use with an off-the-shelf deep-learning-based edge detector [30]. During inference, unlike the training, the input is a sketch. We use the same edge detector to convert it to the edge map for later use. Fig.3 depicts examples of photo, sketches and their corresponding edges. The standardized edge maps have small domain discrepancies. In addition to narrowing the domain gap between the training and test data, the standardization module during inference could narrow the gap of individual sketching styles (e.g., stroke width), which was also similarly shown in [35]. Given that edge maps serve as a proxy for real sketches, we slightly abuse the wording of *synthetic sketches* (or omitted as sketches) hereinafter as they may refer to standardized edge maps.

### 3.2   Reference-Guided Photo Synthesis

Previous works [28, 34] show that photos can be encoded to two disentangled representations: content and style representations. We extend the concept to sketches and show that they can be encoded to disentangled representations. Preserving content representation while replacing the sketch style with a real photo style representation could generate a realistic synthesized photo.

The module is trained in two stages. **1)** Disentangled representation encoding stage learns content and style representations from images via auto-encoding. **2)** We further fine-tune the model with sketch-reference-photo triplets, with regularization loss to guarantee the synthesizing quality. Our model is inspired by
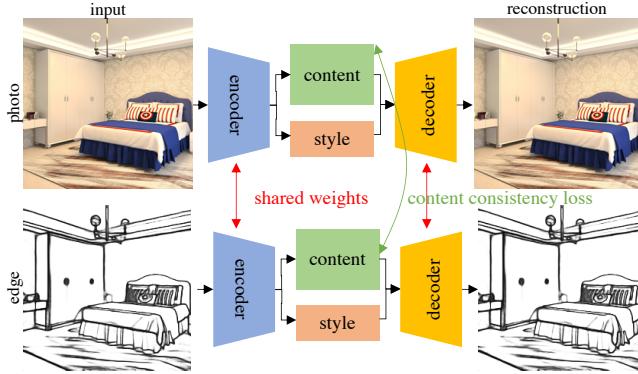
Fig. 4: *Disentangled representation encoding* is the first stage of the sketch-to-photo synthesis module. For each photo, we generate a standardized edge map and form an image pair. Each image of the pair is encoded as content and style representations by the encoder. We add content consistency loss to make content representations of the photo and the edge to be similar. The representations are then decoded to a reconstructed image by the decoder. The network learns the representations through the auto-encoding process. For the performance of sketch to photo synthesis later, both photos and their corresponding standardized edges are fed to the network for auto-encoding.
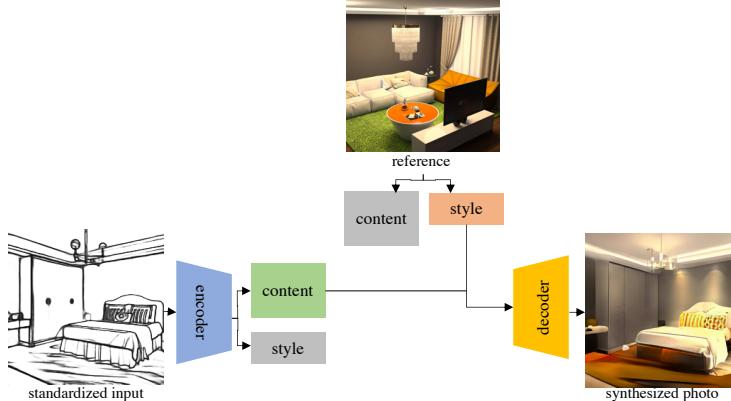


Fig. 5: *Fine-tuning with sketch-reference-photo triplets* is the second stage of the sketch-to-photo synthesis module. The input is a standardized edge map and a reference photo. The model is pre-trained in the representation encoding phase. Both the edge map and the reference photo are encoded by the network for content and style representations. The content and representations are fed to the decoder to reconstruct the synthesized photo.

and based on previous arts on disentangled representation learning [28] and style transfer [34], with novel designs for the goal of scene sketch to photo synthesis. **Disentangled Representation Encoding**. Fig.4 depicts the pipeline of the disentangled representation encoding stage. Denote a pair of input images and its

corresponding edge as $\{\mathbf{x}, \mathbf{x}'\}$, the encoder as $E$, decoder as $G$, and discriminator as $D$. The encoder encodes input pairs $\{\mathbf{x}, \mathbf{x}'\}$ to two representation pairs, content $\{c_{\mathbf{x}}, c_{\mathbf{x}'}\}$ and style $\{s_{\mathbf{x}}, s_{\mathbf{x}'}\}$, i.e., $E(\{\mathbf{x}, \mathbf{x}'\}) = \{\{c_{\mathbf{x}}, c_{\mathbf{x}'}\}, \{s_{\mathbf{x}}, s_{\mathbf{x}'}\}\}$. From the encoded representations, the decoder reconstructs a photo $G(c_{\mathbf{x}}, s_{\mathbf{x}})$ and its edge $G(c_{\mathbf{x}'}, s_{\mathbf{x}'})$. The auto-encoder ensures the reconstructed image pair is similar to the input image pair by the following reconstruction loss in $\ell_1$-norm:

$$\mathcal{L}_{\text{rec}_1} = \mathrm{E}_{\mathbf{x} \sim \mathbf{X}, \mathbf{x}' \sim \mathbf{X}'}[|\mathbf{x} - G(c_{\mathbf{x}}, s_{\mathbf{x}})| + |\mathbf{x}' - G(c_{\mathbf{x}'}, s_{\mathbf{x}'})|] \tag{1}$$

Since the photo and the edge depict the same content, we ask their content representations to be similar in $\ell_1$-norm:

$$\mathcal{L}_{\text{content}} = \mathrm{E}_{\mathbf{x} \sim \mathbf{X}, \mathbf{x}' \sim \mathbf{X}'}[|c_{\mathbf{x}} - c_{\mathbf{x}'}|] \tag{2}$$

Further, the adversarial GAN loss [9] is required to train discriminator $G$ for realistic reconstructions:

$$\mathcal{L}_{\text{GAN}_1} = \mathrm{E}_{\mathbf{x} \sim \mathbf{X}, \mathbf{x}' \sim \mathbf{X}'}[-\log D(G(c_{\mathbf{x}}, s_{\mathbf{x}})) - \log D(G(c_{\mathbf{x}'}, s_{\mathbf{x}'}))] \tag{3}$$

The final loss is $\mathcal{L}_{\text{rec}_1} + \theta \, \mathcal{L}_{\text{content}} + \alpha \, \mathcal{L}_{\text{GAN}_1}$, where $\theta, \alpha$ are both set to be 0.5.

**Fine-Tuning with Sketch-Reference-Photo Triplets**. Fig.5 depicts the pipeline of the fine-tuning stage. Denote the sketch, reference photo and output synthesized photo as $\mathbf{x}^{\mathbf{k}}, \mathbf{x}^{\mathbf{r}}, \mathbf{x}^{\mathbf{o}}$, respectively. With the pre-trained model from the previous representation learning stage, the encoder is able to encode content and style representations of sketches and photos. The output image is generated by the decoder from the content representation of the sketch $c_{\mathbf{x}^{\mathbf{k}}}$, and the style representation of the reference $s_{\mathbf{x}^{\mathbf{r}}}$:

$$\mathbf{x}^{\mathbf{o}} = G(c_{\mathbf{x}^{\mathbf{k}}}, s_{\mathbf{x}^{\mathbf{r}}}) \tag{4}$$

As the model has been pre-trained in the previous stage for encoding content and style representations, the model has a good starting point for synthesizing photos from sketches. To ensure the output image has similar content as the sketch and a similar style as the reference, however, we enforce the following regularization loss on content and style representations in $\ell_1$-norm:

$$\mathcal{L}_{\text{reg}} = \mathrm{E}_{\mathbf{x}^{\mathbf{k}} \sim \mathbf{X}^{\mathbf{k}}, \mathbf{x}^{\mathbf{r}} \sim \mathbf{X}^{\mathbf{r}}, \mathbf{x}^{\mathbf{o}} \sim \mathbf{G}(c_{\mathbf{x}^{\mathbf{k}}}, s_{\mathbf{x}^{\mathbf{r}}})}[|c_{\mathbf{x}^{\mathbf{o}}} - c_{\mathbf{x}^{\mathbf{k}}}| + |s_{\mathbf{x}^{\mathbf{o}}} - s_{\mathbf{x}^{\mathbf{r}}}|] \tag{5}$$

Additionally, the adversarial GAN loss is required:

$$\mathcal{L}_{\text{GAN}_2} = \mathrm{E}_{\mathbf{x}^{\mathbf{k}} \sim \mathbf{X}^{\mathbf{k}}, \mathbf{x}^{\mathbf{r}} \sim \mathbf{X}^{\mathbf{r}}}[-\log D(G(c_{\mathbf{x}^{\mathbf{k}}}, s_{\mathbf{x}^{\mathbf{r}}}))] \tag{6}$$

The final loss is $\mathcal{L}_{\text{reg}} + \beta \, \mathcal{L}_{\text{GAN}_2}$, where $\beta$ is set to be 0.5 in the work.

## 4 Experimental Results

### 4.1 Network Architectures and Training Details

**Network Architectures**. Images are fed to the encoder to obtain content and style representations. First, images go through 4 down-sampling residual blocks
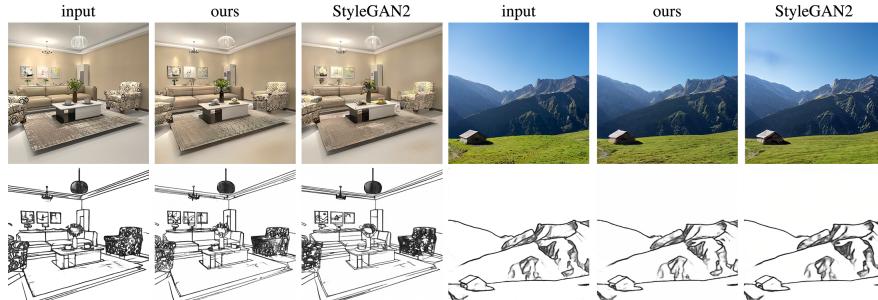
Fig. 6: The reconstruction results of our method and StyleGAN2 [34]. Images are projected into embedding spaces for ours and StyleGAN2 [34]. Both photos and standardized edges are fed to the network for reconstruction. The high faithfulness in reconstruction demonstrates that the learned content and style representations are effective.

Table 1: **(a)** Reconstruction performance measured in LPIPS ($\downarrow$) [41]. Images are projected into embedding spaces for ours and StyleGAN2 [34]. We reconstruct photos and edges with a similar performance as StyleGAN2 [34], demonstrating the disentanglement to content and style representations is effective. **(b)** Reference-guided sketch to photo synthesis performance measured in FID ($\downarrow$) [12]. Our method outperforms other baseline methods in all three categories.

(a)

| input | method | indoor | church | mountain | mean |
|-------|--------|--------|--------|----------|------|
| photo | ours | **0.254** | **0.214** | **0.221** | **0.229** |
| | StyleGAN2 | 0.256 | 0.220 | 0.224 | 0.233 |
| edge | ours | **0.180** | **0.166** | **0.171** | **0.172** |
| | StyleGAN2 | 0.161 | 0.188 | 0.173 | 0.174 |

(b)

| FID ($\downarrow$) | indoor | church | mountain | mean |
|--------------------|--------|--------|----------|------|
| ours | **105.5** | **48.7** | **73.8** | **76.0** |
| SAE [28] | 107.7 | 52.4 | 74.1 | 78.1 |
| ObjSketch [24] | 136.5 | 62.1 | 95.4 | 98.0 |
| SpliceViT [33] | 204.2 | 119.7 | 140.7 | 154.9 |
| DTP [18] | 205.2 | 124.2 | 143.5 | 157.6 |
| Style2Paints [39] | 254.2 | 217.3 | 247.7 | 239.7 |

[11] to obtain an intermediate representation. The intermediate representation is fed to another convolution layer to obtain the content representation with a spatial size of $16 \times 16$. The intermediate representation is also fed to another two convolution layers to obtain a style representation/vector dimension of 2048. The decoder consists of 4 up-sampling residual blocks. The style representation is injected to the decoder convolution layers with weight modulation techniques described in StyleGAN2 [34]. The discriminator is the same as that of StyleGAN2. **Hyper-Parameters and Training Schedules**. For representation encoding, the initial learning rate is 2e-3. We use Adam optimizer [19] with $\beta = (0, 0.99)$. For fine-tuning, we start from the previously pre-trained model. The training schedule stays the same with the initial learning rate being 4e-4. The entire training time for the 3D-front indoor scene dataset is 7 days on 4 V100 GPUs. **Baselines**. We follow the released code and the same settings of all baseline methods and retrain on datasets used in the paper. Specifically, some baselines [18, 24, 28, 33] only work on photos, but not sketches. We use a gray-scale images as a proxy to ensure the photo synthesis quality. Specifically, we first train a

Fig. 7: Various baseline photo syntheses from sketches with style guidance. Note that SpliceViT [33] and DTP [18] are designed for test-time optimization and are not trained on the full dataset, making them disadvantageous to other methods. All other methods are trained on the same dataset with a similar iteration as the proposed method. Style2Paints is designed to synthesize painting, not realistic photos. Our model synthesizes photos that share a similar content as the sketch and a similar visual style as the style photo reference.

sketch to gray-scale photo model using the same setting as step 1 of [24], where the input to the model is a standardized sketch. The generated gray-scale photo is then used to train a gray-scale to color photo model with the same setting of the baseline methods. SpliceViT [33] and DTP [18] are designed for test-time optimization and are not trained on the entire dataset. All other baseline methods are trained on the same dataset as the proposed method with a similar iteration.

### 4.2  Datasets

We train on the following scene photo datasets: **1) 3D-Front Indoor Scene** [6] consists of 14,761 training and 5,479 validation photos. They are rendered with Blender from synthetic indoor scenes including bedrooms and living rooms. Photos are resized to 286 and randomly cropped to 256 during training. **2) LSUN Church** [38] consists of 126,227 photos of outdoor churches. We randomly sample 25,255 photos as the validation set. Photos are resized to 286 and randomly cropped to 256 during training. **3) GeoPose3K Mountain Landscape** [1] has 3,114 mountain landscape photos. 623 photos are randomly sampled for validation. Training photos are resized to 572 and randomly cropped.

For evaluation, we collect a **Scene Sketch Evaluation Set**. For each category (indoor scenes, mountain and church), we collect 50 sketches from the Internet, respectively. The sketches are collected with an intention to cover various sketching styles, e.g. different levels of line width, geometric distortion, use of shading, etc.
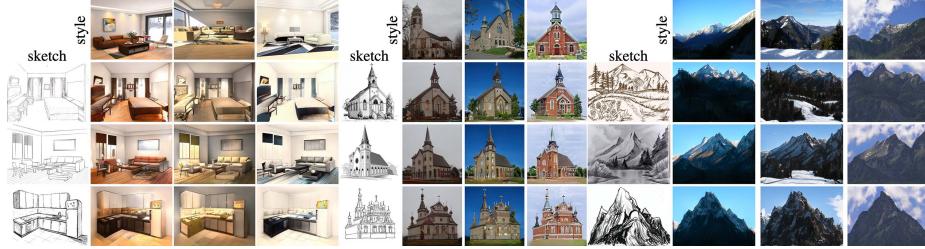
Fig. 8: The indoor scene, church and mountain sketch to photo synthesis with different references. We synthesize high-fidelity scene photos with similar content as the sketch and similar style as the reference photos.

### 4.3   Representation Encoding

With effective learned representation, the model could reconstruct photos or sketches with high quality. We evaluate reconstruction performance in LPIPS [41].

Table 1**a** reports the LPIPS distance of reconstructed and input photos and synthetic sketches of our stage 1 model and StyleGAN2 [34]. Fig.6 depicts several examples of the input and reconstruction. Our representation encoding model has a slightly better reconstruction performance compared to StyleGAN2, indicating the learned content and style representations are adequate and ready for further fine-tuning with sketch-reference-photo pairs.

### 4.4   Photo Synthesis

We evaluate the photo synthesis performance of our method and baselines in terms of photo-realism. We calculate the Fréchet inception distance (FID) [12] between the synthesized photo set and the training photo set for each category (Table 1**b**). Our method outperforms other baselines under the FID metric. Fig.7 depicts synthesis results of our method and baselines. Note that SpliceViT [33] and DTP [18] designed for test-time optimization and was not trained on the full dataset, making it disadvantageous to other methods. Style2Paints is designed to synthesizing painting, not realistic photos. We however include it as it is one of the few works that study synthesizing from scene sketches. Our synthesis result outperforms all other methods, with SAE [28] being the second. As for if the content of the output photo matches with the input sketch or if the style matches with the reference photo, we provide human perceptual evaluation in Section 4.5.

We also provide more visualization of our synthesis results of indoor scenes, churches and mountains in Fig.8.

### 4.5   Human Perceptual Study

We conduct a human perceptual study to evaluate the realism of synthesized photos, and if synthesized photos match contents and styles as desired. We only

Table 2: A human perceptual study of the synthesized photos. **(a)** The fooling rate of our synthesized model over real photos measures the realism of the generation. **(b)** User preference on which method synthesizes photos that depicts more similar content to the sketch. **(c)** User preference on which method synthesizes photos that depicts more similar visual style to the reference photo. Compared with [28], we have a higher fooling rate over real photos, better content and style matching preference rate.

(a) Fooling rate (↑)      (b) Content matching (↑)      (c) Style matching (↑)

| (%) | indoor scene | church | mountain | mean |
|---|---|---|---|---|
| ours | **25.00** | **44.3** | **48.9** | **39.4** |
| SAE [28] | 10.0 | 6.6 | 20.0 | 12.2 |

| (%) | indoor scene | church | mountain | mean |
|---|---|---|---|---|
| ours | **80.1** | **92.1** | **75.0** | **82.4** |
| SAE [28] | 19.9 | 7.9 | 25.0 | 17.6 |

| (%) | indoor scene | church | mountain | mean |
|---|---|---|---|---|
| ours | **61.9** | **90.9** | **71.0** | **74.6** |
| SAE [28] | 38.1 | 9.1 | 29.0 | 25.4 |

evaluate our method and SAE [28], the second best-performing synthesis method, due to limited resources.

We create a survey consisting of three parts: photorealism, content matching with sketches and style matching with reference photos. As guidance to the participants, we state our research purpose at the beginning of the survey. For each part, a detailed description and an example question with answers and explanations are provided for the participant's reference. The order of our results, baseline results, and real images are randomly shuffled in the survey to minimize the potential bias from the participant. Each part consists of 13 questions, with one question being a *bait question* with an obvious answer. The bait question is designed to check if the participant is paying attention and if the answers are reliable. There are in total 51 participants, with 1 being ruled out due to failing one of the bait questions. Thus we finally collect 1,950 valid human judgments.

To evaluate the photorealism, we randomly select synthesized photos of ours and SAE evenly from three categories. Both methods use the same input sketch and reference photo. For each synthesized photo, we use Google's search by image feature to find the most similar real photo and ask participants which one they think looks more like a real photo. We then calculate the percentage of participants being fooled. Note that the fooling rate of random guessing is 50%. Table **2a** reports the fooling rate of our method and SAE. Ours is 27% higher than SAE. Specifically, for churches and mountains, ours achieves a fooling rate over 44%: the generated photos are almost indistinguishable from real photos.

To evaluate if the synthesized photos match the content of the input sketch, we show participants an input sketch and two synthesized results from our method and SAE, and ask them to pick one that has the most similar content as the sketch. Table **2b** reports the preference rate of ours over SAE. We achieve 82% on average preference rate, well outperforming the baseline.

To evaluate if the synthesized photos match the style of the reference photo, we show participants a reference photo and two synthesized results from our method and SAE, and ask them to pick one that has the most similar style to the sketch. Table **2c** reports the preference rate of ours over SAE. We achieve a 75% average preference rate, well outperforming the baseline.

Fig. 9: The style representations of sketches and photos are well separated, while the content representations of sketches and photos are tangled together. We visualize learned content and style representations of sketches and photos with T-SNE [25]. The results show that sketches and photos share the content space and it is appropriate to train on photos and transfer knowledge to sketches.
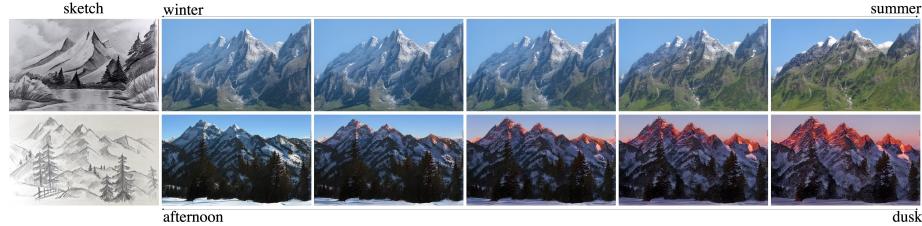


Fig. 10: Sketch to photo synthesis with combined style representations of two references. We encode style representations from two photos, e.g. a winter photo and a summer photo. By increasing the weight of the summer image and decreasing that of the winter image, the synthesized photo from the sketch gradually changes from winter appearance to summer appearance.

### 4.6    Photo Editing Through Sketch

As depicted in Fig.11, given an input photo, we convert it to a standardized edge map (where we refer as sketch for simplicity). Users could add and remove strokes to edit the photo. We also show the possibility of sequential editing in the figure. We evaluate the photo editing performance for the indoor scene validation dataset, and the FID [12] of edited images to the training set is 69.2. One limitation is that the content in the unmodified region of a given photo may not be well preserved as the edited photo is solely generated from the edge map.

### 4.7    Analysis and Ablation Studies

**Analysis of Style Representations**. We visualize the learned content and style representations of photos and sketches using T-SNE [25] in Fig.9: style representations of sketches and photos are well separated, while content representations of sketches and photos are not separable. This verifies the grounding of the method: the content representations of sketches and photos can be shared, while the style representations for the two are different. Thus, combining the
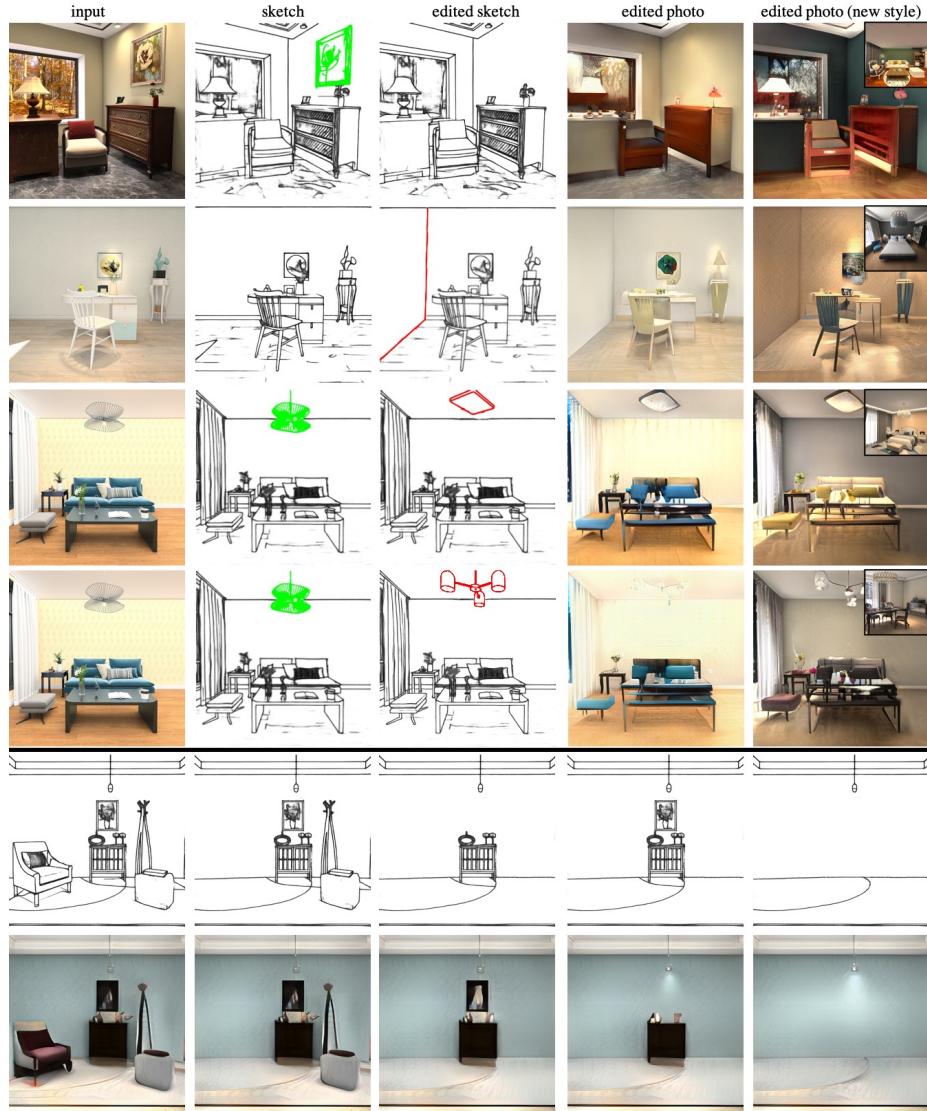
| input | sketch | edited sketch | edited photo | edited photo (new style) |

Fig. 11: Photo editing and style transfer via sketches. *Upper:* Given an input image, we first convert it to a standardized edge map. We then add or remove strokes in the edge map and convert it back to a photo. The visual style of the photo could also be changed with a reference photo (top right). *Lower:* Sequential editing by gradually removing strokes.

content representation of a sketch and style representation of a photo could decode a realistic synthesized photo.

**Style Interpolation**.We study if the reference style can be a combination of style of two different reference images $\mathbf{x^{r_1}}$ and $\mathbf{x^{r_2}}$. Suppose their style representations

Table 3: Ablation studies on the fine-tuning stage, content and style regularization loss for indoor scenes in FID ($\downarrow$) [12] distance. Having both stage 2 fine-tuning and the regularization loss gives the best result.

| no fine-tune | fine-tune+style loss | fine-tune+content loss | fine-tune+all loss |
|---|---|---|---|
| 107.9 | 107.0 | 106.1 | **105.5** |

are $\mathbf{s_{x^{r_1}}}$ and $\mathbf{s_{x^{r_2}}}$. The combined representation $s_{\text{combined}} = \gamma \mathbf{s_{x^{r_1}}} + (1-\gamma)\mathbf{s_{x^{r_2}}}$, where $\gamma \in [0,1]$. By adjusting $\gamma$, we synthesize photos with a combined style from both reference images. Fig.10 depicts examples of mountain sketch to photo synthesis with combined styles from two different reference images. By adjusting $\gamma$, the synthesized photos have a continuous interpolation from winter to summer, and afternoon to dusk.

**Fine-Tuning Model**. One of the novelty is that we propose the fine-tuning with sketch-reference-photo triplets for the task. We evaluate if the fine-tuning is necessary by removing the fine-tuning stage. As reported in Table 3, removing the model fine-tuning leads to 2.4 worse results in the FID metric.

**Content and Style Regularization Loss**. We study if the regularization loss at the fine-tuning stage is effective. We study the function of the content loss ($|c_{\mathbf{x^o}} - c_{\mathbf{x^k}}|$) and style loss ($|s_{\mathbf{x^o}} - s_{\mathbf{x^r}}|$) respectively. As reported in Table 3, removing the content regularization loss leads to 1.5 worse results in FID metric, and removing the style loss leads to 0.6 worse results. This verifies the effectiveness of the proposed regularization loss.

## 5  Summary

We propose a reference-guided framework for photo synthesis from scene sketches. We first convert all input photos and sketches to standardized edge maps, allowing the model to learn in unsupervised setting without the need of real sketches or sketch-photo pairs. Sequentially, the standardized input and reference image are disentangled into content and style components to synthesize new hybrid image that preserves the content of standardized input while transferring the style of reference image. Extensive experiments demonstrate that our method can generate and edit a realistic photo from a user's scene sketch with a reference photo as style guidance, surpassing the previous approaches on three benchmarks.

A major insight of this work is that, we learn to synthesize scene structures directly from the vast amount of readily-available photos, rather than synthesizing and combining individual objects. Rather than worrying about the acclimated errors from sketch-based object detection, photo synthesis and spatial combination for the final output, we treat the scene sketches as a whole and learn the holistic structures for photo synthesis.

One limitation is that the deep-learning based standardization step could eliminate strokes that reflect the details of the scene, or misinterpret the strokes as textures. Future work could study a sketch-to-edge standardization process that preserves higher fidelity of the sketch. Another limitation lies in the sketch-based photo editing - the unchanged regions of a given photo may not be well

preserved. This is due to the model takes sketch as the only input. Future work could improve the performance by taking the original photo into consideration.

# References

1. Brejcha, J., Čadík, M.: Geopose3k: Mountain landscape dataset for camera pose estimation in outdoor environments. Image and Vision Computing **66**, 1–14 (2017)
2. Chen, S.Y., Su, W., Gao, L., Xia, S., Fu, H.: Deepfacedrawing: Deep generation of face images from sketches. ACM Transactions on Graphics (TOG) **39**(4), 72–1 (2020)
3. Chen, W., Hays, J.: Sketchygan: Towards diverse and realistic sketch to image synthesis. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 9416–9425 (2018)
4. Choi, Y., Uh, Y., Yoo, J., Ha, J.W.: Stargan v2: Diverse image synthesis for multiple domains. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 8188–8197 (2020)
5. Delanoy, J., Aubry, M., Isola, P., Efros, A.A., Bousseau, A.: 3d sketching using multi-view deep volumetric prediction. Proceedings of the ACM on Computer Graphics and Interactive Techniques **1**(1), 1–22 (2018)
6. Fu, H., Cai, B., Gao, L., Zhang, L.X., Wang, J., Li, C., Zeng, Q., Sun, C., Jia, R., Zhao, B., et al.: 3d-front: 3d furnished rooms with layouts and semantics. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10933–10942 (2021)
7. Gao, C., Liu, Q., Xu, Q., Wang, L., Liu, J., Zou, C.: Sketchycoco: Image generation from freehand scene sketches. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 5174–5183 (2020)
8. Ghosh, A., Zhang, R., Dokania, P.K., Wang, O., Efros, A.A., Torr, P.H., Shechtman, E.: Interactive sketch & fill: Multiclass sketch-to-image translation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1171–1180 (2019)
9. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. Advances in neural information processing systems **27** (2014)
10. Guérin, É., Digne, J., Galin, E., Peytavie, A., Wolf, C., Benes, B., Martinez, B.: Interactive example-based terrain authoring with conditional generative adversarial networks. Acm Transactions on Graphics (TOG) **36**(6), 1–13 (2017)
11. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
12. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in neural information processing systems **30** (2017)
13. Huang, X., Liu, M.Y., Belongie, S., Kautz, J.: Multimodal unsupervised image-to-image translation. In: Proceedings of the European conference on computer vision (ECCV). pp. 172–189 (2018)
14. Iizuka, S., Simo-Serra, E., Ishikawa, H.: Globally and locally consistent image completion. ACM Transactions on Graphics (ToG) **36**(4), 1–14 (2017)
15. Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1125–1134 (2017)
16. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 4401–4410 (2019)

17. Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and improving the image quality of stylegan. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 8110–8119 (2020)
18. Kim, S., Kim, S., Kim, S.: Deep translation prior: Test-time training for photorealistic style transfer. arXiv preprint arXiv:2112.06150 (2021)
19. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: ICLR (Poster) (2015)
20. Larsson, G., Maire, M., Shakhnarovich, G.: Learning representations for automatic colorization. In: European conference on computer vision. pp. 577–593. Springer (2016)
21. Lee, H.Y., Tseng, H.Y., Mao, Q., Huang, J.B., Lu, Y.D., Singh, M., Yang, M.H.: Drit++: Diverse image-to-image translation via disentangled representations. International Journal of Computer Vision **128**(10), 2402–2417 (2020)
22. Ling, H., Kreis, K., Li, D., Kim, S.W., Torralba, A., Fidler, S.: Editgan: High-precision semantic image editing. Advances in Neural Information Processing Systems **34** (2021)
23. Liu, M.Y., Huang, X., Mallya, A., Karras, T., Aila, T., Lehtinen, J., Kautz, J.: Few-shot unsupervised image-to-image translation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10551–10560 (2019)
24. Liu, R., Yu, Q., Yu, S.X.: Unsupervised sketch to photo synthesis. In: European Conference on Computer Vision. pp. 36–52. Springer (2020)
25. Van der Maaten, L., Hinton, G.: Visualizing data using t-sne. Journal of machine learning research **9**(11) (2008)
26. Meng, C., He, Y., Song, Y., Song, J., Wu, J., Zhu, J.Y., Ermon, S.: Sdedit: Guided image synthesis and editing with stochastic differential equations. In: International Conference on Learning Representations (2021)
27. Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y.: Semantic image synthesis with spatially-adaptive normalization. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 2337–2346 (2019)
28. Park, T., Zhu, J.Y., Wang, O., Lu, J., Shechtman, E., Efros, A., Zhang, R.: Swapping autoencoder for deep image manipulation. Advances in Neural Information Processing Systems **33**, 7198–7211 (2020)
29. Pathak, D., Krahenbuhl, P., Donahue, J., Darrell, T., Efros, A.A.: Context encoders: Feature learning by inpainting. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2536–2544 (2016)
30. Poma, X.S., Riba, E., Sappa, A.: Dense extreme inception network: Towards a robust cnn model for edge detection. In: Proceedings of the IEEE/CVF winter conference on applications of computer vision. pp. 1923–1932 (2020)
31. Portilla, J., Simoncelli, E.P.: A parametric texture model based on joint statistics of complex wavelet coefficients. International journal of computer vision **40**(1), 49–70 (2000)
32. Richardson, E., Alaluf, Y., Patashnik, O., Nitzan, Y., Azar, Y., Shapiro, S., Cohen-Or, D.: Encoding in style: a stylegan encoder for image-to-image translation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2287–2296 (2021)
33. Tumanyan, N., Bar-Tal, O., Bagon, S., Dekel, T.: Splicing vit features for semantic appearance transfer. arXiv preprint arXiv:2201.00424 (2022)
34. Viazovetskyi, Y., Ivashkin, V., Kashin, E.: Stylegan2 distillation for feed-forward image manipulation. In: European Conference on Computer Vision. pp. 170–186. Springer (2020)

35. Wang, J., Lin, J., Yu, Q., Liu, R., Chen, Y., Yu, S.X.: 3d shape reconstruction from free-hand sketches. arXiv preprint arXiv:2006.09694 (2020)
36. Wang, L., Qian, C., Wang, J., Fang, Y.: Unsupervised learning of 3d model reconstruction from hand-drawn sketches. In: Proceedings of the 26th ACM international conference on Multimedia. pp. 1820–1828 (2018)
37. Xiang, X., Liu, D., Yang, X., Zhu, Y., Shen, X., Allebach, J.P.: Adversarial open domain adaptation for sketch-to-photo synthesis. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 1434–1444 (2022)
38. Yu, F., Seff, A., Zhang, Y., Song, S., Funkhouser, T., Xiao, J.: Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. arXiv preprint arXiv:1506.03365 (2015)
39. Zhang, L., Li, C., Simo-Serra, E., Ji, Y., Wong, T.T., Liu, C.: User-guided line art flat filling with split filling mechanism. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2021)
40. Zhang, R., Isola, P., Efros, A.A.: Colorful image colorization. In: European conference on computer vision. pp. 649–666. Springer (2016)
41. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 586–595 (2018)
42. Zhou, Y., Zhu, Z., Bai, X., Lischinski, D., Cohen-Or, D., Huang, H.: Non-stationary texture synthesis by adversarial expansion. arXiv preprint arXiv:1805.04487 (2018)
43. Zhu, J.Y., Krähenbühl, P., Shechtman, E., Efros, A.A.: Generative visual manipulation on the natural image manifold. In: European conference on computer vision. pp. 597–613. Springer (2016)
44. Zhu, S.C., Wu, Y., Mumford, D.: Filters, random fields and maximum entropy (frame): Towards a unified theory for texture modeling. International Journal of Computer Vision **27**(2), 107–126 (1998)
45. Zou, C., Yu, Q., Du, R., Mo, H., Song, Y.Z., Xiang, T., Gao, C., Chen, B., Zhang, H.: Sketchyscene: Richly-annotated scene sketches. In: Proceedings of the european conference on computer vision (ECCV). pp. 421–436 (2018)