

# FaceSketching: Interactive Realistic Face Image Creation from Free-hand Sketches

Anonymous Author(s)  
Submission Id: 1234

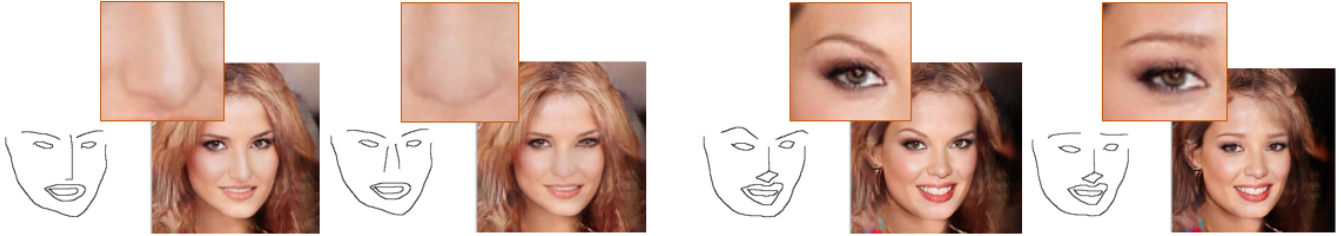


Figure 1: This is a teaser

## ABSTRACT

### CCS CONCEPTS

• **Computer systems organization** → **Embedded systems**; *Redundancy*; Robotics; • **Networks** → Network reliability.

### KEYWORDS

Image synthesis, sketch-based interface, face editing, deep neural network

#### ACM Reference Format:

Anonymous Author(s). 2020. FaceSketching: Interactive Realistic Face Image Creation from Free-hand Sketches. *ACM Trans. Graph.* 1, 1 (May 2020), 6 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

Flexibly creating new content is one of the most important goals in both computer graphics and computer-human interaction. While sketching is an efficient and natural way for common users to express their ideas for designing and editing new content, sketch-based interaction techniques have been extensively studied since the very early stage of computer graphics [Chen et al. [n.d.], 2008; Igarashi et al. 1999; Sutherland 1964; Zeleznik et al. 1996]. Imagery content is the most ubiquitous media with a large variety of display devices everywhere in our daily life. Creating new imagery content is one way to show people's creativity and communicate smart ideas. In this paper, we target portrait imagery, which is inextricably bound to our life, and present a sketch-based system that allows common users to create new face imagery and edit them by specifying the desired facial shapes via free-hand strokes.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2020 Association for Computing Machinery.  
0730-0301/2020/5-ART \$15.00  
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

Deep learning techniques bring significant improvement on the realism of virtual images. Recently, a large number of studies have been conducted on face image editing. (Add papers.) All these methods require a real image as the start points. When creating a non-existent face, sketching is the most natural way to describe the desired structure and shapes. Treating free-hand sketches as a painting style, image translation techniques can be directly applied. (Cite image translation and style transfer papers.) However, the challenges in image translation is the edge alignment problem. The ambiguities and large discrepancy from real face shapes existed in imperfect free-hand sketches poses significant challenges for image translation networks which have strong power on realistic texture synthesis but perform poorly on geometric shape creation. Moreover, local editing is not supported in these uninterpretable end-to-end networks.

(Space and scale varying ambiguity in sketches) In order to support user's control on the creation and local editing from imperfect sketches, the ambiguities of hand-drawn strokes which vary in different facial regions at different scales should be solved for generating geometrically correct guidance. (Spatially-varying.) A user may draw the face contour in a casual way or depict the eye or eyebrow shape carefully in little ambiguity. (Temporally-varying.) As immediately see the generated face image, the user may add or edit local strokes with more careful control. Therefore, the ambiguity and distortion in different stage of the sketching process vary from the beginning to the end.

(Is speed an issue for sketch-base interface?)

(Contributions) In summary, our contribution in this paper is three-fold.

- Based on comprehensive analysis on the edge alignment issue in image translation DNNs, we propose an sketch-to-face system that allows users to control the desired facial shapes in adaptive precision.
- A spatial varying feature normalization in deep neural network to support global creation and local editing in multiple scales.

- Robust to various styles of hand-drawn strokes and allows progressively refinement to generate desired shapes and facial details.
- An easy-to-use system for face creation and editing which requires hand-drawn sketches only as input without any other specification of region masks.

## 2 RELATED WORK

### 2.1 Image-to-Image Translation

Given an input image from one domain, an image-to-image translation model outputs a corresponding image from another domain and preserves the content in the input image. Existing image-to-image translation models are based on generative adversarial networks conditioned on images. Pix2pix [?] is the first general image-to-image translation model which is able to be applied to different scenarios according to the paired training images, such as, semantic maps to real images, day images to night images, image coloring, and edge maps to real images. [?] utilizes semantic label maps and attributes of outdoor scenes as input and generates the corresponding photo-realistic images. In order to model multi-modal distribution of output images, BicycleGAN [?] encourages the connection between the output and the latent code to be invertible. CycleGAN [?], DualGAN [?], and DiscoGAN [?] propose unsupervised image translation model with a similar idea named cycle consistency borrowed from language translation literature. StarGAN [?] presents a framework for one-to-many image translation by adding a domain code as input and a domain classifier as guidance. Pix2pixHD [?] is proposed as a high-resolution image-to-image translation model for generating photo-realistic image from semantic label maps using a coarse-to-fine generator and a multi-scale discriminator. It can also be applied to edge-to-photo generation by using the paired edge maps and photos as training data. However, xxxxxxxxxxxxxxx

### 2.2 Sketch-based Image generation

Sketch-based image generation is a hot topic in computer vision and computer graphics. Given a sketch of a scene with text labels for objects, traditional methods, such as Sketch2Photo [?] and PhotoSketcher [?], search corresponding image patches from a large image dataset and then fuse the the retrieved image patches together according to the sketch. These methods are not able to ensure the global consistency of the resultant image and fails to generate totally new images. After the breakthrough made by deep neural networks (DNN) in computer graphics and computer vision, a variety of DNN-based models have been proposed for sketch-based image generation. The general image-to-image translation models mentioned above are able to be applied to sketch-based image generation once sketches and the corresponding images are used as training data. Besides, many other models are designed specially for sketch inputs. SketchyGAN [?] aims to generate real images from multi-class sketches. A novel neural network module, called mask residual unit (MRU), is proposed to improve the information flow by injecting the input image at multiple scales. Edge maps are extracted from real images and tiled as training sketches. However, the resultant images of SketchyGAN are still not satisfied.

### 2.3 pooling

### 2.4 Face Generation and Editing

All existing methods for local face editing require the user to provide masks and strokes manually. In comparison, strokes are only the input to indicate the desired shape. Our system automatically interprets the intended edit and produces the local change accurately. This greatly reduces users' burden and preserves the fluency of user interaction.

Only local edit in a relatively small area is supported. As reported in SC-FEGAN [], artifacts appear when complete a large region in FaceShop [].

## 3 DEEP NETWORK FOR SKETCH-PHOTO TRANSLATION

### 3.1 Edge alignment in baseline Model

Generating a realistic photo from sketch can be considered as an image translation problem. We use the state-of-the-art Pix2PixHd [] as our baseline.

We use the CelebA-HD dataset [] which contains xxx face images in WxH. All the face images are globally aligned according to their eye positions. For each photo, we generate sketch samples by xxx method. XX for training and xx for testing. Using this paired sketch-photo dataset, we trained our method and other state-of-the-art approaches for comparison.

**3.1.1 Global alignment.** While all the face images in the CelebA dataset [] are globally aligned with their eye positions, the learned generator implicitly embeds the global layout. Once the drawn sketches deviate from this implicitly embedded layout, the learned translation model generate awkward results, (as Fig. 3 shows).

**3.1.2 Data augmentation with geometric translation.** A straightforward way is to augment the training set by random geometric transformation of input sketches. However, large interval/range of the transform yields un-convergence of the network training. We limit the transformation range to  $(-\frac{\pi}{10}, \frac{\pi}{10})$  rotation,  $(-\frac{W}{20}, \frac{W}{20})$  translation, and (scale?) in order to increase the tolerance of the trained model on the distortion and roughness of hand-drawn sketches. Moreover, as an interactive system, we also simply provide a reference sketch, like ShadowDrawing []. We place an averaged face contour image on the canvas to provide a reference coordinate system for user to draw their strokes. Therefore, the drawn sketches under this geometric reference will be located in or close to space of the training sketches.

**3.1.3 Face Sketches and Data Augmentation.** Paired face sketch-photo dataset is required for supervised sketch-to-face generation methods. Since there exists no large-scale face sketch datasets, the training face sketches used by existing methods are generated from face image dataset, e.g. CelebA-HQ face dataset, using edge detection algorithm such as HED [?]. However, the sketches generated by edge detection algorithm are sometimes incomplete or other problems. (I would say the edge maps are quite different from handdrawn sketches, not because of the incompleteness.) Although CSAGAN [?] applied self-attention mechanism to alleviate the incompleteness problem, the others remains. Therefore we use another method to generate clearer and complete sketches from

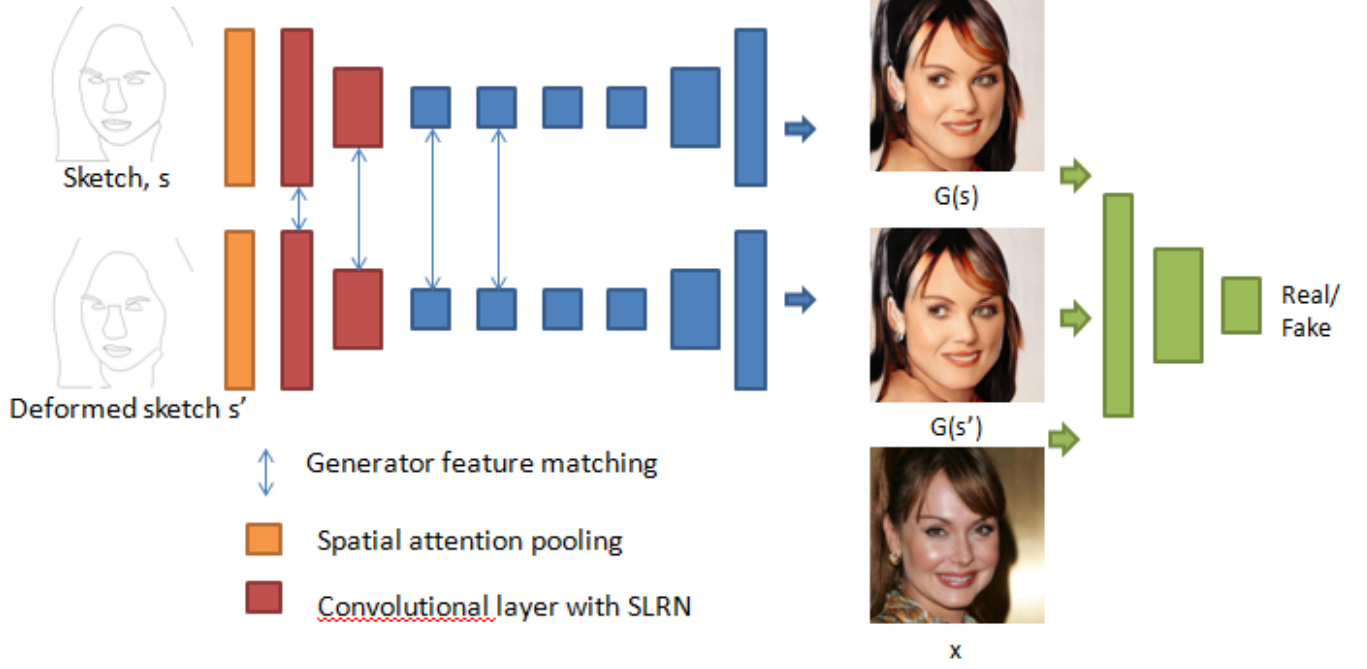


Figure 2: The architecture of our model.

Figure 3: Generated face images from a sketch that does not follow the globally aligned face layout.

face images. The CelebAMask-HQ dataset [?] provides a face semantic map for each face image in CelebA-HQ dataset. We basically use the boundary map of the semantic map as the sketch of the corresponding face image. Figure 4 shows an example of comparison between a sketch generated by edge detector and from semantic boundary. A contour2photo model is trained in Pix2PixHD to generate face images from face contours that contains facial makers only. Hair shapes, or other facial features such as beard, wrinkles, are not supported. (Compare with contour data? With different styles of input sketches, the trained model is not well generalized. )

**Stroke Deformation.** A shortcut of sketches generated from semantic boundary (and those generated by edge detector) is the lines of sketches are perfectly aligned to the edges of the corresponding face images. In order to break the edge-alignment and mimic the hand-drawn sketches, we apply a deform to the lines. Specifically, we vectorize the lines of each sketches using Auto-Trace algorithm [?], and add an offset randomly selected from  $[-d, d]^2$  to the control points and end points of the vectorized lines, where  $d$  is the maximum offset and we set  $d = 11$  in our experiments.

Figure 4: Comparison between a sketch generated from edge detection and from semantic boundary.

**3.1.4 pix2pixHD without instance normalization.** The baseline model, as well as many existing stylization DNNs, uses spatial normalization (instance normalization?) to extract style statistics, treating them as spatially uniform on the entire image. The instance normalization is defined as

$$f^i(x, y, c) = \frac{f^i(x, y, c)}{\sum_{x, y} f^i(x, y)}, \quad (1)$$

where  $i, x, y, c$  respectively indicates feature map layer, pixel position, and channel indexes. However, the shallow convolution layers typically extract texture or brightness statistics, which are varying dramatically on different local regions of a hand-drawn sketch due to its sparsity. For example, there might be a large number of strokes around hair regions or mouth regions to describe details. These heavy strokes bring significant difference while instance normalization at each convolution layer normalizes local patch features with a global factor. Therefore, the extracted features for identical eye strokes could be very different. This instance normalization term leads to significant change of local shapes, as shown in Fig. ??.(add a figure here or show comparison in result section. )

We think that strokes mainly describe the shape features, without little texture information. Normalization is expected to remove

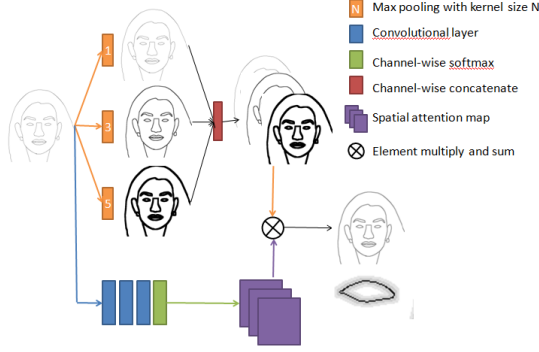


Figure 5: Sap

the variance of global stroke styles, such as stroke widths, stroke curveness, and stroke coarseness. During painting, the stroke coarseness is spatially varying when user depict different facial features. The user may draw very carefully for features such as eyes, eyebrows, but not so accurate on cheek shapes.

**3.1.5 Spatial Attention Pooling.** When the input hand-drawn sketch is not well-drawn, it is a trade-off between the realism of the output face image and the alignment between the input sketch and the output face image. In order to alleviate the edge alignment between the input sketch and the output face image, we should relax the sharp sketch lines with one-pixel width. (relax thin strokes to an ambiguity band with various width or uncertainty.) One of the straightforward ways is to smooth the lines of sketches using image smoothing algorithm. Another is to dilate the sketch lines so that the widths of lines are of multiple pixels [?]. (if it is not officially published, it is not necessary to cite it. but if our method performs better, then cite it and show comparison.) However, the capacity of either the two hand-crafted ways above is limited, because the uniform smoothness and the dilate radius for all positions of the whole sketch violate the unevenness of hand-drawn sketches on depicting different facial parts. We argue that the balance between the realism of the output face image and the alignment between the input sketch and the output face image differs from one position to another across the face image. Therefore, the smoothness or the dilation radius should be spatial-specific.

Based on the discussion above, we propose a new module, called spatial attention pooling (SAP), to dilate the input sketch in a spacial-specific way. (I would not use 'dilate' since this simple word does not reveal the underlying discovery.) Given an input sketch  $S \in \mathbb{R}^{H \times W}$ , we first pass it through  $N_r$  pooling branches with different kernel sizes of  $r_i, i = 1, \dots, N_r$  to get  $\{P_i | i = 1, \dots, N_r\}$ . Then we compute the spatial attention map  $W \in \mathbb{R}^{N_r \times H \times W}$  with  $W = \text{Softmax}(f(s))$ , where  $f()$  is implemented with two convolutional layers. ( $f(S)$  to extract low-level features from the input sketch?) (It is not clear what is the motivation of computing the SA map  $W$ . It might be better to describe your idea of how to use  $W$  with  $P$  and then describe how to get  $W$ .) A softmax layer which is computed over channels is added at the end of the convolutional layers, ensuring that for each position, the sum of weights of all channels equals to 1. The output of SAP is computed as:



Figure 6: Spatial attention pooling to balance edge alignment and stroke ambiguity at different facial regions.

$$SAP(s) = \sum_{i=1}^{N_r} W_i * P_{r_i}(s), \quad (2)$$

where  $W_i$  is the  $i$ th channel of  $W$ , and  $*$  is element-wise multiplication. We show this idea in Fig. 6.

### 3.1.6 Losses. generator feature matching effect and losses summary

(Before describe how you do it, please describe why you do this.)

Let  $G_q(\cdot)$  produces the feature maps of the  $q$ -th layer in the generator  $G$ . Given an input sketch  $S$  and the corresponding deformed sketch  $\tilde{S}$ , we compute the generator feature matching loss as:

$$\mathcal{L}_{GFM}(G) = \mathbb{E}_{S \sim p_{data}(S)} \frac{1}{N_Q} \sum_{q \in Q} \frac{1}{|G_q|} \|G_q(S) - G_q(\tilde{S})\|_1, \quad (3)$$

where  $|G_q|$  denotes the number of elements in  $G_q(\cdot)$ ,  $Q$  indicates a set of the selected generator layers for computing this loss and the size of  $Q$  is  $N_Q$ . We select the xxx layers of the generator in our experiments.

Besides the generator feature matching loss  $\mathcal{L}_{GFM}(G)$ , for generator  $G$  and multi-scale discriminator  $D = D_k | k = 1, 2, \dots, N_D$ , the adversarial loss  $\mathcal{L}_{GAN}(G, D)$  and the discriminator feature matching loss  $\mathcal{L}_{DFM}(G, D)$  are computed as the same form as those in pix2pixHD [?]. Discussion: add equations of these two losses or not The objective of the proposed model is:

$$\min_G \max_D \mathcal{L}_{GAN}(G, D) + \lambda \mathcal{L}_{DFM}(G, D) + \mu \mathcal{L}_{GFM}(G). \quad (4)$$

where  $\lambda$  and  $\mu$  are the weights for balancing different losses. We set  $\lambda = \text{xxx}$  and  $\mu = \text{xxx}$  in our experiments.

## 4 EXPERIMENTS

### 4.1 Sketch Interface

We develop a sketch-based interface which allows common users to describe their desired face and part shapes by a few strokes. Once the user finishes his drawing, the generated face image is shown on the right bottom of the sketch? This will be helpful for local editing., as Fig. 7 shows. The user is allowed to edit the sketch by erasing strokes or adding new stroke to change eye shapes, noses shapes, eyebrows, and so on. Each round of face image generation





Figure 7: Sketch interface. Full interface with editing tools.



Figure 8: Face generation with different models. From left to right: (a) hand-drawn sketches as input. (b) Results generated by pix2pixHD that is retrained using our contour-photo dataset. (c) Results generated by pix2pixHD that is retrained using our contour-photo dataset with geometric transformation as data augmentation. (d) Results generated by removing instance normalization at the shallow convolution layers (five layers in the global generator). (e) Results from our model (pix2pixHD architecture by replacing instance normalization with the proposed SLRN.) More results can be found here: (Provide a link for all results.)

after users' modification takes about ... seconds on a XXX GPU with xxGB memory.

## 4.2 Face Generation from Contours

*Contour Dataset.* (Yangbinxin: Add description on the dataset generation.)

*Photo Generation from Contours.* The Pix2pixHD [ ] was applied to the task of translating edges to photos. They successfully generate high-quality photos from contours that are generated from real photos. However, when we apply the Pix2pixHD model with hand-drawn sketches which present distinct characteristics from the synthesized contours which are smooth and clean without geometric distortions, it fails to produce good results, as Fig. 8 shows. (Add more results with other models.) We can see that (expected results: (e) have better details. (d) i have no idea. c is better than (b).)

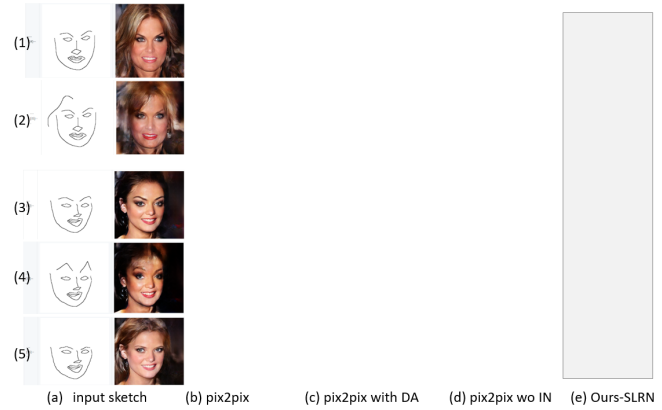


Figure 9: Local face editing with different models. The proposed SLRN captures the shape details in the drawn sketches and successfully avoid edge-aligned artifacts caused by distortion in hand-drawn sketches. More results can be found here: (Provide a link for all results.)

*Face Editing with Strokes.* When users want to change local shapes of facial features, a few strokes can be modified in our interface. The pix2pixHD model does not take the sparsity of sketches and the instance normalization tends to normalize local regions with a global style factor, which mainly conveys brightness variance. Therefore, when local details such as hair, the bamoustache, and so on, the results change significantly even only local region is modified. In comparison, our proposed SLRN captures the shape details in the drawn sketches and successfully avoid edge-aligned artifacts caused by distortion in hand-drawn sketches, as Fig. 9 shows.

In order to analyze the effect of instance normalization and our SLRN, we visualize the features extracted at early stages in the generators of different models. As Fig. 10 shows, ...

## 4.3 Comparison with Image Translation networks

Existing DNNs for image translation can be trained for sketch-photo translation using the paired dataset.

## 4.4 Comparison with Image Editing

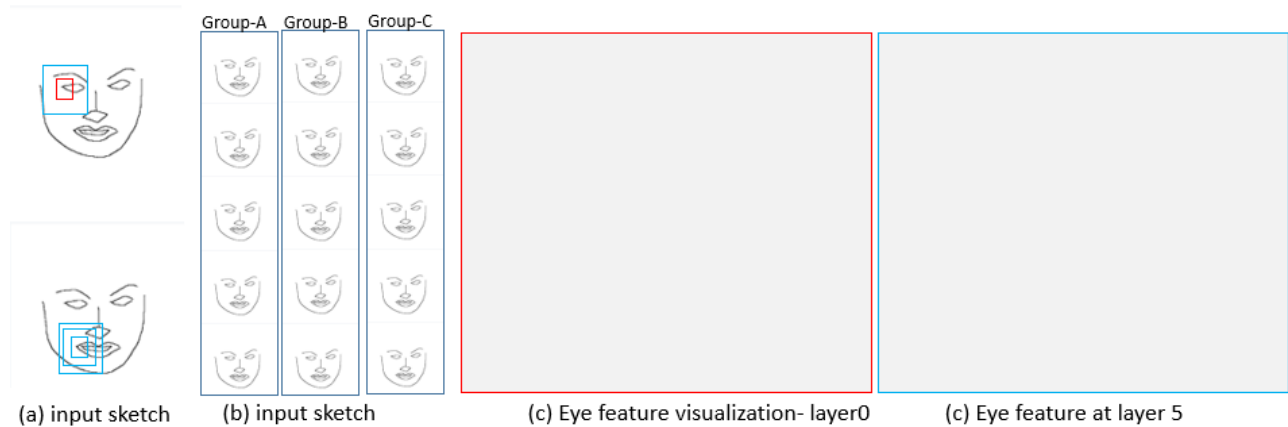
## 4.5 Limitations and future work

Add color

Combine with attribute

## REFERENCES

- Tao Chen, Ming ming Cheng, Ping Tan, Ariel Shamir, and Shi min Hu. [n.d.]. Sketch2Photo: Internet image montage. *ACM Trans. Graph* ([n. d.]).
- Xuejin Chen, Sing Bing Kang, Ying qing Xu, and Julie Dorsey. 2008. Sketching reality: Realistic interpretation of architectural designs. *ACM Trans. Graph.* 27, 2 (2008), 11:1–11:15.
- Takeo Igarashi, Satoshi Matsuoka, and Hidehiko Tanaka. 1999. Teddy: A sketching interface for 3-D freeform design. In *Proc. SIGGRAPH*. 409–416.
- Ivan E. Sutherland. 1964. Sketch pad a man-machine graphical communication system. In *DAC*.
- Robert C. Zeleznik, Kenneth P. Herndon, and John F. Hughes. 1996. SKETCH: An Interface for Sketching 3D Scenes. In *Proc. SIGGRAPH*. 163–170.



**Figure 10: Visualization of extracted features under local editing.** We extract the features at different convolution layers at the left eye position (a) from three groups of sketches (b) with different types of local editing. The extracted high-dimensional features using different models including pix2pixHD-DA, pix2pixHD-wo-IN, and our SLRN) are mapped into 2D space using TSNE [?] in (c). (ChengZhiHua: provide a link for all results.)