

000  
001  
002054  
055  
056

# Face Sketch Synthesis with Style Transfer using Pyramid Column Feature

003  
004  
005  
006  
007  
008  
009  
010  
011057  
058  
059  
060  
061  
062  
063  
064  
065

Anonymous ICCV submission

Paper ID 612

012  
013066  
067

## Abstract

014  
015  
016  
017  
018068  
069  
070  
071

In this paper, we propose a novel framework based on deep neural networks for face sketch synthesis from a photo. Imitating the process of how artists draw sketches, our framework synthesizes face sketches in a cascaded manner. A content image is first generated that outlines the shape of the face and the key facial features. Textures and shadings are then added to enrich the details of the sketch. We utilize a fully convolutional neural network (FCNN) to create the content image, and propose a style transfer approach to introduce textures and shadings based on a newly proposed pyramid column feature. We demonstrate that our style transfer approach based on the pyramid column feature can not only preserve more sketch details than the common style transfer method, but also surpasses traditional patch based methods. Quantitative and qualitative evaluations suggest that our framework outperforms other state-of-the-arts methods, and can also generalize well to different test images.

019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031072  
073

## 1. Introduction

032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093

Face sketch synthesis has drawn a great attention from the community in recent years because of its wide range of applications. For instance, it can be exploited in law enforcement for identifying suspects from a mug shot database consisting of both photos and sketches. Besides, face sketches have also been widely used for entertainment purpose. For example, filmmakers could employ face sketch synthesis technique to ease the cartoon production process.

045  
046  
047  
048  
049  
050  
051  
052  
053094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

Unfortunately, there exists no easy solution to face sketch synthesis due to the big stylistic gap between photos and sketches. In the past two decades, a number of exemplar based methods [14, 11, 19, 20] were proposed. In these methods, a test photo is first divided into patches. For each test patch, a candidate sketch patch is identified by finding the most similar photo patch in a training set of

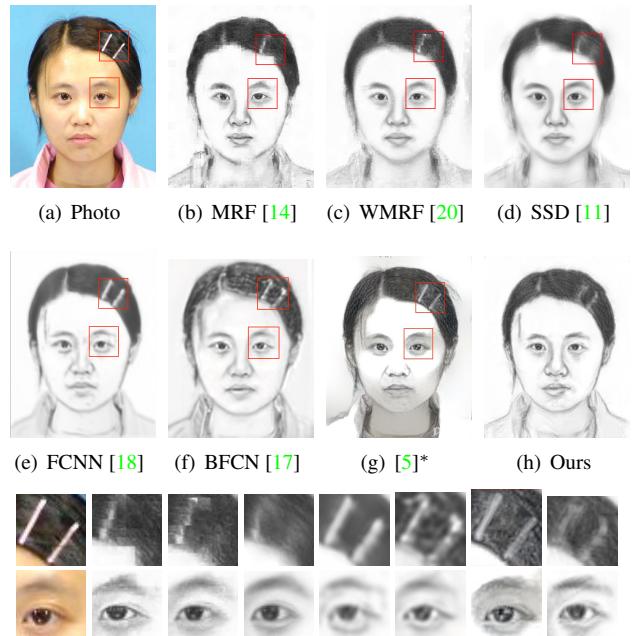


Figure 1. Face sketches generated by existing methods and the proposed method. Our method can not only preserve both hair and facial content, but also contains sharp textures. (Note that (g) is obtained from the deep art website<sup>1</sup> using the photo as content and a sketch from the training set as style.)

photo-sketch pairs. The main drawback of this approach is that if there exists no photo patch in the training set which is sufficiently similar to a test patch, loss of content will be observed in the synthesized sketch. For example, the sketches in the first row of Fig. 1 fail to keep the hairpins. Besides, some methods [11, 20] blur away the textures when they try to eliminate the inconsistency between neighboring patches. Another common problem is that the synthesized sketch may not look like the test photo (see the left eye in Fig. 1(b)). Recently, approaches [17, 18] based on convolutional neural network (CNN) were developed to solve these problems. Since they directly generate sketches from photos, structures and contents of the photos can be maintained. However, the pixel-wise loss functions adopted

<sup>1</sup><https://deeprart.io/>

108 by these methods will lead to blurry artifacts (see Fig. 1(e)  
109 and 1(f)) because they are incapable of preserving texture  
110 structures. The popular neural style transfer provides a better  
111 solution for texture synthesis. However, there exist two  
112 obstacles in directly applying such a technique. First, the  
113 result is easily influenced by the illumination of the photo  
114 (see the face in Fig. 1(g)). Second, it requires a style image  
115 to provide the global statistics of the textures. If the given  
116 style image does not match with the target sketch (which we  
117 do not have), some side effects will occur (see the nose in  
118 Fig. 1(g)). (*ken: I cannot see the problem in 1(h)*)

119 For an artist, the process of sketching a face usually starts  
120 with outlining the shape of the face and the key facial features  
121 like the nose, eyes and mouth. Textures and shadings  
122 are then added to regions such as hair, lips, and bridge of the  
123 nose to give the sketch a specific style. Based on the above  
124 observation, and inspired by neural style transfer [4], we  
125 propose a new framework for face sketch synthesis from a  
126 photo that overcomes the aforementioned limitations. In our  
127 method, a content image that outlines the face is generated  
128 by a feed-forward neural network, and textures and shadings  
129 are then added using a style transfer approach. Specifically,  
130 we design a new architecture of fully convolutional  
131 neural network (FCNN) composed of inception layers [12]  
132 and convolution layers with batch normalization [6] to  
133 generate the content image (see Section 4.1). To synthesize  
134 the textures, we first divide the target sketch into a grid.  
135 For each grid cell, we compute a newly proposed pyramid  
136 column feature using the training set (see Section 4.2). A  
137 target style can then be computed from a grid of these pyramid  
138 column features, and applied to the content image. Our  
139 approach is superior to the current state-of-the-art methods in  
140 that

- 141 • It is capable of generating more stylistic sketches without  
142 introducing over smoothing artifacts.
- 143 • It can preserve the content of the test photo well.

## 144 2. Related Work

### 145 2.1. Face Sketch Synthesis

146 Based on the taxonomy of previous studies [11, 20],  
147 face sketch synthesis methods can be roughly categorized  
148 into profile sketch synthesis methods [1, 2, 15] and shading  
149 sketch synthesis methods [8, 11, 13, 14, 18, 19, 20].  
150 Compared with profile sketches, shading sketches are more  
151 expressive and thus more preferable in practice. Based on  
152 the assumption that there exists a linear transformation  
153 between a face photo and a face sketch, the method in [13]  
154 computes a global eigen-transformation for synthesizing a  
155 face sketch from a photo. This assumption, however, does  
156 not always hold since the modality of face photos and that  
157 of face sketches are quite different. Liu et al. [8] pointed out

158 that the linear transformation holds better locally, and therefore  
159 they proposed a patch based method to perform sketch  
160 synthesis. In [14], a MRF based method was proposed to  
161 preserve large scale structures across sketch patches. Variants  
162 of the MRF based method were introduced in [19, 20]  
163 to improve the robustness to lighting and pose, and to render  
164 the ability of generating new sketch patches. In addition  
165 to these MRF based methods, approaches based on guided  
166 image filtering [11] and feed-forward convolutional neural  
167 network [18] are also found to be effective in transferring  
168 photos into sketches. A very recent work similar to ours is  
169 reported by Zhang et al. [17]. They proposed a two-branch  
170 FCNN to learn content and texture respectively, and then  
171 fused them through a face probability map. Although their  
172 results are impressive, their sketch textures do no look natural  
173 and the facial components are over smoothed.

### 174 2.2. Style Transfer with CNN

175 Texture synthesis has long been a challenging task. Traditional  
176 methods can only imitate repetitive patterns. Recently,  
177 Gatys et al. [4, 5] studied the use of CNN in style  
178 representation, and proposed a method for transferring the  
179 style of one image (referred to as the style image) to another  
180 (referred to as the content image). In their method, a target  
181 style is first computed based on features extracted from  
182 the style image using the VGG-Network. An output image  
183 is then generated by iteratively updating the content image  
184 and minimizing the difference between its style and the target  
185 style. Justin et al. [7] further accelerated this process by  
186 learning a feed forward CNN in the training stage. These  
187 methods represent styles by a multi-scale Gram matrix of  
188 the feature maps. Since the Gram matrix only cares about  
189 global statistics, local structures may be destroyed when  
190 the style image is very different from the content image.  
191 Although this may not be a problem in transferring artistic  
192 styles to images, this will definitely produce noticeable  
193 artifacts in face sketch synthesis as people are very sensitive  
194 to the distortions of the facial features. In [3], Chen and  
195 Schmidt proposed a different patch based style transfer  
196 method which is better at capturing local structures. However,  
197 it is still far from satisfactory to be employed in face  
198 sketch synthesis. Our style transfer approach is inspired by  
199 but different from the above work [4, 5, 7] in that our target  
200 style is computed from image patches of many different  
201 images rather than just from one single image.

## 202 3. Style Representation

203 Following the work of [5], we use Gram matrices of  
204 VGG-19 [10] feature maps as our style representation. Denote  
205 the vectorized  $c$ th channel of the feature map in the  $l$ th  
206 layer of the final sketch  $\mathcal{X}$  by  $F_c^l(\mathcal{X})$ . A Gram matrix for  
207 the feature map in the  $l$ th layer is then defined by the inner

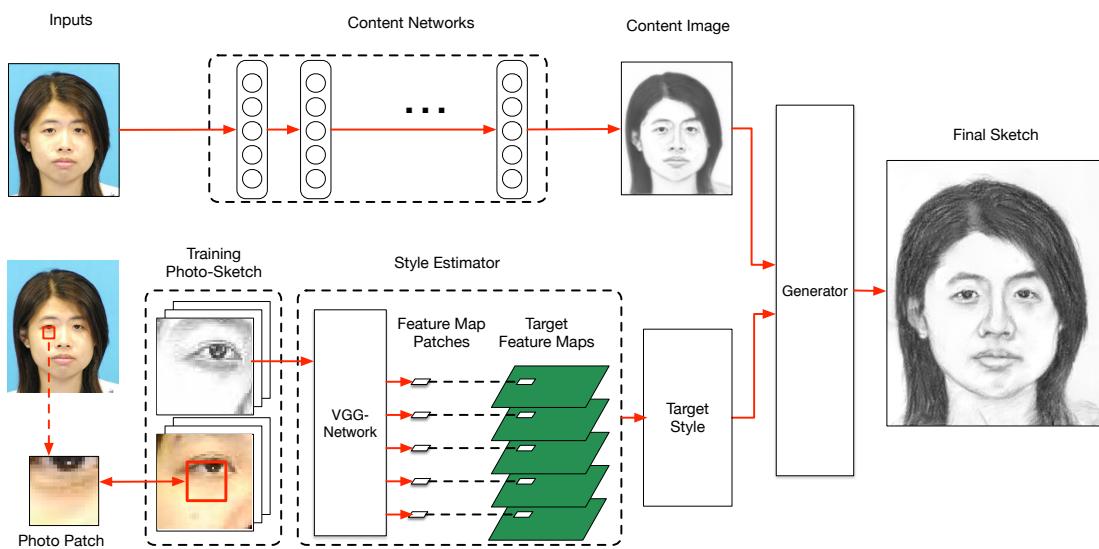


Figure 2. The proposed method contains two branches which take an eye-aligned test photo as input. The content network outputs a content image which outlines the face, and the style estimator generates a target style. The final sketch is generated by combining the target style with the content image.

products between two channels of this feature map, i.e.,

$$G_{ij}^l(\mathcal{X}) = F_i^l(\mathcal{X}) \cdot F_j^l(\mathcal{X}), \quad (1)$$

where  $G^l(\mathcal{X}) \in \mathcal{R}^{N_l \times N_l}$  and  $N_l$  is the number of channels of the feature map in the  $l$ th layer. Since  $G_{ij}^l(\mathcal{X})$  is an inner product between two channels of the feature map, a Gram matrix is actually a summary statistics of the feature map without the spatial information. Generally, a Gram matrix of the feature map captures the density distribution of a sketch. For example, if a given style (sketch) image has much less hair than the test photo, the synthesized sketch  $\mathcal{X}$  will become brighter than a natural sketch (see experimental results in Section 5.1). Thus it is important to have a style (sketch) image which is (statistically) similar to the test photo. Note that, in face sketch synthesis, however, there usually does not exist a single photo-sketch pair in the training set that matches all properties of the test photo. How to compute a target style for the synthesized sketch  $\mathcal{X}$  is therefore not trivial, and is the key to the success of this approach. We will introduce a feature-space patch-based approach to solve this problem in Section 4.2.

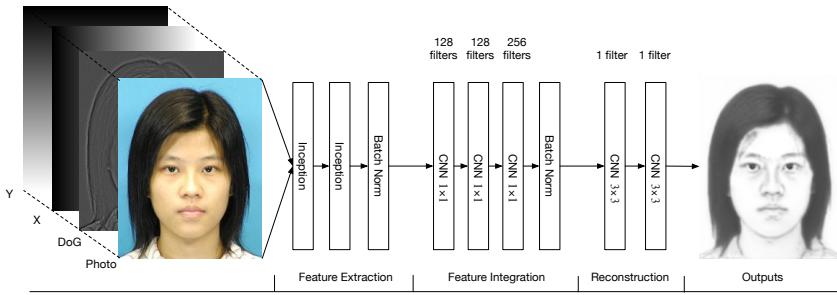
## 4. Methodology

Our method can be classified as a shading synthesis method. The steps of our method are summarized in Fig. 2. First, a preprocessing step as described in [14] is carried out to align all photos and sketches in the training set by the centers of the two eyes. An eyes-aligned test photo  $\mathcal{I}$  is then fed into two branches, namely the content network and the style generator. The content network converts  $\mathcal{I}$  into a content

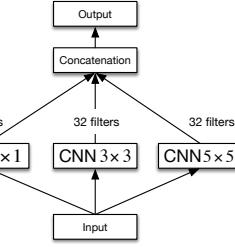
image  $\mathcal{C}$ , which outlines the shape of the face and the key facial features such as nose, eyes, mouth and hair. The style estimator divides  $\mathcal{I}$  into non-overlapping  $16 \times 16$  patches. For each test patch, it locates the most similar photo patch from the photo-sketch pairs in the training set and produce a target sketch patch from the corresponding sketch in the pair. A pyramid column feature (Section 4.2) is then computed for the target sketch patch. Finally, a target style can be computed from a grid of these pyramid column features, and the final sketch can be synthesized by applying the target style to  $\mathcal{C}$  through neural style transfer [5].

### 4.1. Content Image Generation

Our content network architecture is shown in Fig. 3. Together with the test photo, we feed two extra channels containing spatial information (i.e., x and y coordinates) and a difference of Gaussian (DoG) image into the content network. As pointed out in [14], face sketch synthesis algorithms benefit from integrating features from multiple resolutions. We employ an inception module inspired by the GoogLeNet [12] to extract features. It concatenates feature maps generated from filters with different spatial resolutions. Our inception unit contains 3 different size of filters ( $1 \times 1$ ), ( $3 \times 3$ ) and ( $5 \times 5$ ) (see Fig. 3(b)). Then, the output features are fed to a three-layer-CNN for feature integration, where the size of all filters are fixed at  $1 \times 1$ . Finally, the integrated features are used to reconstruct the content map by a two-layer-CNN with the filter size being  $3 \times 3$ . Because  $L_1$ -norm is better at preserving details than  $L_2$ -norm, we use the  $L_1$ -norm between the generated content image  $\mathcal{C}$



(a) The architecture of content network



(b) Inception module

Figure 3. Illustration of the content network for generating a content image. The numbers above the building block denote the number of CNN filters. (a) The architecture of content network. (b) The inception module in (a) contains three groups of filters with different sizes.

and ground truth sketch  $S$  as the loss function, see Equ. 2.

$$\mathcal{L}_s = \frac{1}{N} \sum_{i=1}^N \|C_i - S_i\| \quad (2)$$

where  $N$  is the number of photos.

## 4.2. Pyramid Feature Column

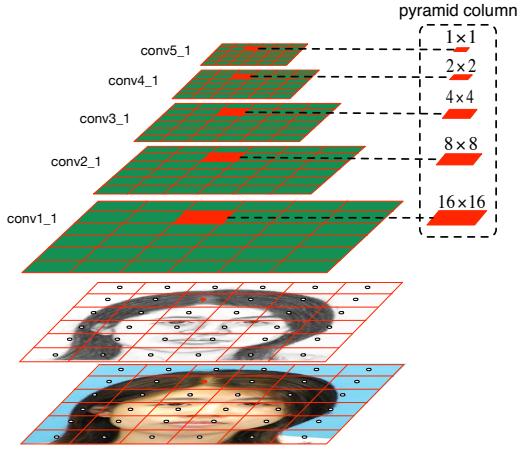
Fig. 4 shows an example of the pyramid feature column. Denote the feature maps (of the  $l$ th layer) used to estimate the style of the final sketch by  $A^l$ . In our feature patch based method, we divide  $A^l$  into a fixed size of grid. Due to the different feature map size, the sizes of the feature patches at layer  $conv1\_1$ ,  $conv2\_1$ ,  $conv3\_1$ ,  $conv4\_1$  and  $conv5\_1$  are  $16 \times 16$ ,  $8 \times 8$ ,  $4 \times 4$ ,  $2 \times 2$  and  $1 \times 1$ . The photos and sketches are resized to  $288 \times 288$ , thus the size of grid is  $18 \times 18$ . Grouping feature map patches having the same grid indexes  $(i, j)$  at different layers together, we get a pyramid feature column  $U_{ij}$ . To estimate  $U_{ij}$ , a sketch patch in the training set is fed to the VGG-Network and a pyramid column is composed of the resulting feature maps. This process consists of two steps: (1) find a matching sketch patch  $S_{ij}$  from the training set for  $U_{ij}$  and (2) feed  $S_{ij}$  to VGG-Network and extract  $U_{ij}$  from the resulting feature maps.

**Find Matching Sketch Patch** Similar to previous works [14, 20], we examine the similarity of test photo patch and train photo patches to find a matching sketch patch  $S_{ij}$  for  $U_{ij}$ . Denote the train photo patch set as  $P$ , given a test photo patch  $T'_{ij}$ , the target sketch patch  $S_{ij}$  is found by

$$S_{ij} = \phi(\min_{T_{mn}^k \in P} f(T_{mn}^k, T'_{ij})) \quad (3)$$

$$m = i + \Delta x, n = j + \Delta y \quad (4)$$

where  $T_{mn}^k$  is the photo patch of the  $k$ th image at  $(m, n)$  in training set,  $\phi$  is a one-to-one mapping between photo patch and sketch patch,  $\Delta x$  and  $\Delta y$  is the shift of  $(m, n)$  around  $(i, j)$ . The function  $f$  measures the discrepancy between

Figure 4. Illustration of pyramid feature column. Feature maps of the final sketch  $A^l$  are divided into a fixed  $18 \times 18$  grid. A pyramid column  $U_{ij}$  consists of feature map patches at different layer having the same grid indexes  $(i, j)$ .

$T_{mn}^k$  and  $T'_{ij}$ . Because we only need the test photo patch and train photo patch have a similar appearance, we simply use MSE as  $f$  here.

**Estimate Pyramid Feature Column** After we get  $S_{ij}$ , we can compute the corresponding pyramid feature column  $U_{ij}$ . Although the feature column can be generated by a  $16 \times 16$  patch, padding zeros is not a good idea. Therefore, we will make use of the surrounding area of  $S_{ij}$  to compute  $U_{ij}$ . According to the architecture of VGG-19, the receptive field of  $conv5\_1$  is 132 without padding. To get  $16 \times 16$  patches, we will need a  $144 \times 144$  region  $\mathcal{R}$  containing 81 patches in total.  $S_{ij}$  should be at the center of  $\mathcal{R}$ , i.e.  $S_{55}$  with respect to the grids of  $\mathcal{R}$ . For those border part where we can't find a region centered at  $S_{ij}$ , we just find another region  $\mathcal{R}'$  which has the biggest intersection with the should be region  $\mathcal{R}$  (the blue box in Fig. 5). The  $S_{pq}$  in  $\mathcal{R}$  should be  $S'_{pq}$  in  $\mathcal{R}'$  ( $(p, q)$ ) is calculated according to different con-

ditions). The feature column  $U_{55}$  in  $\mathcal{R}$  or  $U_{pq}$  in  $\mathcal{R}'$  is the estimated pyramid feature column.

Allocating all these feature columns to its original patch location  $(i, j)$ , we can get the final feature maps  $\tilde{U}$  and use it to compute gram matrix.

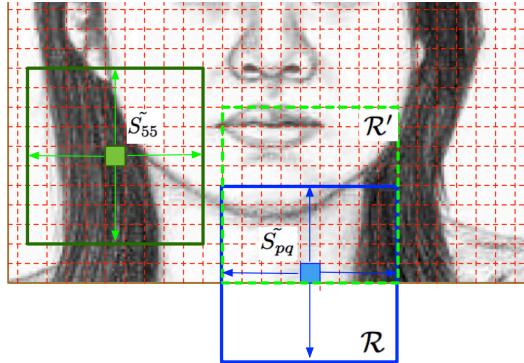


Figure 5. Find a region  $\mathcal{R}$  centered at  $S_{ij}$ . If  $\mathcal{R}$  is inside the image,  $S_{ij}$  should be the center patch of  $\mathcal{R}$ , see the left region. Otherwise, we need to find  $\mathcal{R}'$  inside the image which has the biggest intersection with  $\mathcal{R}$ , see the right part. Suppose the indexes of  $S_{ij}$  inside  $\mathcal{R}'$  is  $(p, q)$ , then  $U_{pq}$  is the corresponding feature column.

### 4.3. Loss Function of Sketch

The same as [5], our loss function has a style loss and content loss. In addition we add a component loss to enhance the key facial components. The total loss is

$$\mathcal{L}_t(\mathcal{X}) = \alpha \mathcal{L}_c + \beta_1 \mathcal{L}_s + \beta_2 \mathcal{L}_k \quad (5)$$

We minimize the loss function by updating the target sketch  $\mathcal{X}$  without changing the VGG parameters.

The content loss is defined based on the difference between the feature map of the sketch and the content image at layer conv1\_1:

$$\mathcal{L}_c(\mathcal{X}) = \|F^{\text{conv1\_1}}(\mathcal{X}) - F^{\text{conv1\_1}}(\mathcal{C})\|_2^2. \quad (6)$$

And the style loss is the difference between the gram matrix of the final sketch and the target gram matrix

$$\mathcal{L}_s(\mathcal{X}) = \sum_{l \in L_s} \frac{1}{M_l^2 N_l^2} \|G^l(\mathcal{X}) - G^l(\tilde{U})\|_2^2 \quad (7)$$

where  $N_l$  denotes the number feature maps at layer  $l$ , and  $M_l$  is the feature map width times height.

To better transfer styles of the key facial components, we employ a component loss to encourage the key component style of the final sketch being the same as the target key component style. Since two eyes are placed at fixed positions, the key components lie roughly within a rectangular region taking the positions of two eyes as vertices. Key component style is given by gram matrices calculated

within feature map regions  $\mathcal{K}$  corresponding to the key components. These regions are specified by pyramid feature columns  $U_{ij}$  whose  $C_{ij}$  are inside  $\mathcal{K}$ . We take out  $\mathcal{K}$  and calculate the corresponding gram matrix loss as the component loss

$$\mathcal{L}_k(\mathcal{X}) = \sum_{l \in L_s} \frac{1}{\hat{M}_l^2 N_l^2} \|\hat{G}^l(\mathcal{K}) - \hat{G}^l(\tilde{U}_k)\|_2^2 \quad (8)$$

where  $\hat{M}_l$  denotes height times width of  $\mathcal{K}$ .

### 4.4. Implementation Details

**VGG-19 Parameters** Since the VGG-Network is originally designed for color images, while sketches are gray scale images, we modify the first layer of VGG-Network for gray scale images by setting the filter weights to

$$W^k = W_r^k + W_g^k + W_b^k \quad (9)$$

where  $W_r^k$ ,  $W_g^k$ , and  $W_b^k$  are weights of the  $k$ th filter in the first convolutional layer for the R, G and B channels respectively, and  $W^k$  is the weight of the  $k$ th filter in the first convolutional layer of our modified network.

**Data Partition** CUHK has 88 training photos and 100 test photos, and AR has 123 photos. Our training set is composed of the 88 training photos of CUHK and 100 photos from AR. When training the content network, 10% of the training set are taken out as validation set. All the 188 photo-sketch pairs are used to generate target sketch.

**Training content network** The input photo-sketch pairs are aligned by the facial landmarks of eyes, and all images are resized to  $200 \times 256$ . A mirror padding is carried out before the convolution operation when necessary to ensure the output feature map is the same size as the input. Adadelta [16] is used as the optimizer because it is stable and much faster than others.

**Sketch Generation** In all experiments, we resize test photos and photo-sketch pairs in the training set to a fixed size of  $288 \times 288$ . The final sketch is obtained by resizing the resulting sketch back to the original size. The size of  $\mathcal{K}$  is  $48 \times 48$ ,  $\alpha = 0.004$ ,  $\beta_1 = 1$  and  $\beta_2 = 0.1$ . The minimization is carried out using L-BFGS. Instead of using random noises, we use the content image as a starting point, which will make the optimization process converge much faster.

## 5. Experiments

We evaluate the performance of the proposed method against other state-of-the-art methods on the CUHK student

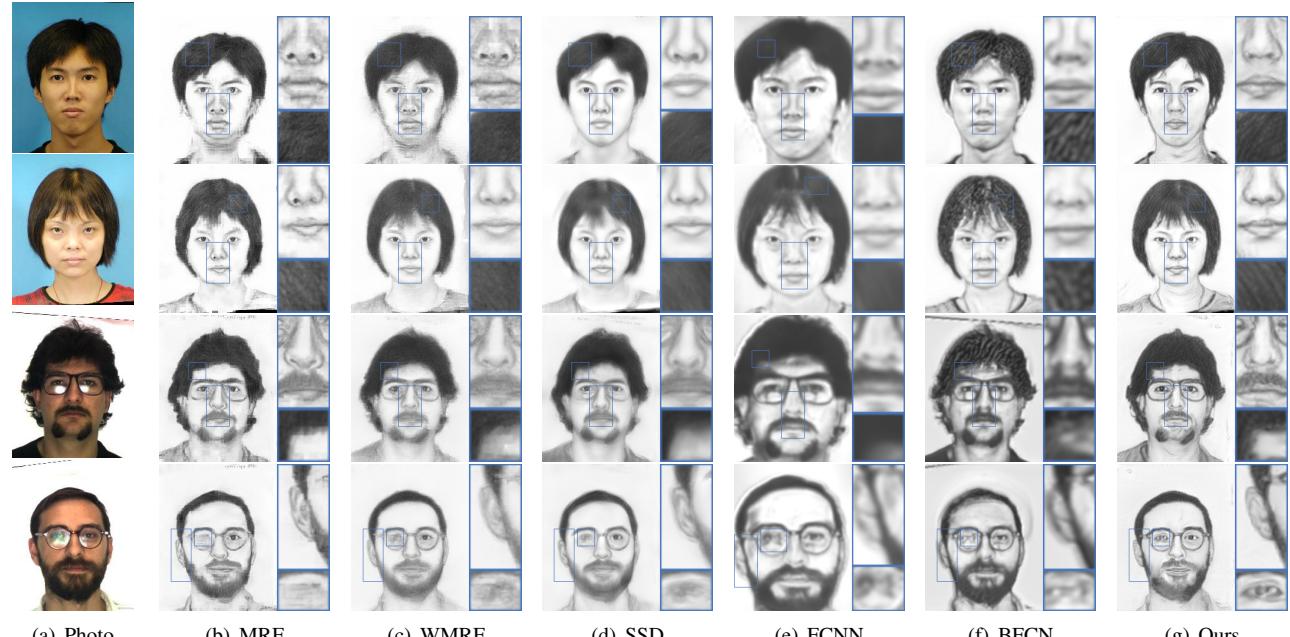


Figure 6. Examples of qualitative evaluation on CUHK (first two rows) and AR (last two rows). (a) The sketches drawn by artists. (b) MRF [14] (c) WMRF [20] (d) SSD [11] (e) FCNN [18] (f) BFCN [17] (g) Ours. The proposed method preserves more texture details, for example in the hair and nose. It is also best at keeping the origin structures of photos, such as the glasses.

dataset [14] and the AR dataset [9]. We compare the results of our method against other 6 methods including traditional approaches and recent deep learning models. After discussing the disadvantages of previous quantitative evaluation criteria, we introduce the Normalized Gram Matrix Difference (NGMD) as a new evaluation tool. We believe that NGMD is more effective and will greatly promote the sketch qualities in future work.

## 5.1. Style Transfer Evaluation

Although style transfer has shown remarkable performance in artistic style, it can't be directly applied to face sketch synthesis, see Fig. 1(h). The generated sketch is greatly influenced by illumination of original photo and doesn't look like sketch at all. To prove our assumption in Section 3 that gram matrix captures the density distribution of sketch, we replace the original RGB photo with content image generated by our content network. We selected 3 sketch styles with different amount of hairs and see how it influences the result. Fig. 7 shows a clear relationship between the hair amount of sketch and pixel intensities of sketch result. The facial key parts in Fig. 7(c) and 7(d) are missed. Therefore, even with a content structure, the original style transfer is still not suitable for elaborate task. However, if the style image has a similar structure with test photo, for example Fig. 7(b), the results can be quite good. This inspires our feature patch based method. Image patch method is not considered because it will introduce patch in-

consistency as discussed in Section 2.

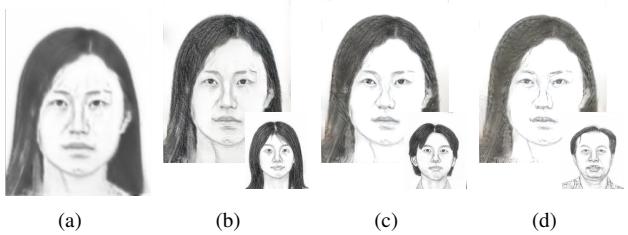


Figure 7. (a) is the content image generated by our content network. (b), (c) and (d) are generated by different styles. It can be seen that when the hair decreases the generated sketch becomes brighter.

## 5.2. Sketch Generation

Fig. 6 shows the comparison between our methods and many other approaches. The first two rows are from CUHK test set, and the last two are from AR. We can see that our method can generate more stylistic sketches than others. For example, in the hair part, only MRF, BFCN and the proposed method can generate obvious textures. However, the texture of MRF is not continuous in the border part and introduce many other artifacts, and the texture of BFCN doesn't look like human strokes. Both WMRF and SSD introduce over smoothing effect and FCNN is not able to give clear textures. Our method can not only generate textures for hairs and mustache but also shadings, for example the

648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661

nose part.

On the other hand, only BFCN and the proposed method can handle structures decorated on the face well, for example the glasses of the last two row. MRF, WMRF and SSD are exemplar based method so they can't handle something different from training set, the glass edges of them are not complete. FCNN, BFCN and our method generate the image content by CNN, so they can handle the original photo structures well. But both FCNN and BFCN care nothing about the textures of facial part, so their results are not good considering the face part. In contrast, our method can well maintain the image content and attach textures similar to human paintings.

### 662 5.3. Quantitative Results

Methods	AR			CUHK		
	R1	R5	R10	R1	R5	R10
FCNN	-	-	-	81%	96%	97%
MRF	97.5%	97.5%	100%	83%	96%	96%
WMRF	97.5%	97.5%	100%	83%	97%	98%
SSD	96.7%	97.5%	100%	87%	97%	98%
Ours	98.4%	98.4%	100%	87%	98%	99%

671 Table 1. Recognition rate on benchmark datasets. The best performance is colored in red.

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

**Sketch Recognition** Sketch synthesis methods are usually evaluated quantitatively via the face sketch recognition task [11, 14, 18, 20]. If an algorithm achieves higher sketch recognition rates, it suggests that this method is more effective in synthesizing sketches. We adopt the widely used PCA based recognition method with “rank-1 (R1)”, “rank-5 (R5)” and “rank-10 (R10)” criteria [14] where “rank  $n$ ” measures the rate of the correct answer in the top  $n$  best matches. The results of different methods are shown in Table 1. Our method achieves the best performance against all other methods in the “R1” and “R5” tests.

However, such kind of evaluation is not effective since every one is close to 100%, and more importantly, it can't guide our research effort. Zhang *et al.* [18] proposed a Multiscale Pixel-wise Reconstruction Loss (MPRL). Although we agree that if MPRL is zero, the generated sketches will be exactly the same with ground truth, there are two reasons why it is not a good criteria. First, experiment results of [18, 17] show that  $L_2$ -norm would give blurry sketches. Second, the ground truth draw by human usually doesn't correspond exactly with photos, thus difficult for algorithms to reach the ground truth.

Meanwhile, human can often be able to tell whether a sketch is generated by algorithms or drawn by artists just from a quick glance without examining details. This indicates that there exists a big gap between the style of algorithm-generated sketches and that of those drawn by

artists. And we need a criteria which can quantitatively measure such gap.

**Normalized Gram Matrix Difference (NGMD)** To quantitatively evaluate the style similarity between a generated sketch and the sketch drawn by a real artist, we employ the normalized gram matrix difference (NGMD) between the generated sketch  $\mathcal{X}$  and the drawn sketch  $\tilde{\mathcal{X}}$ :

$$D_s^l(\mathcal{X}, \tilde{\mathcal{X}}) = \frac{\|G^l(\mathcal{X}) - G^l(\tilde{\mathcal{X}})\|_2^2}{\|G^l(\tilde{\mathcal{X}})\|_2^2} \quad (10)$$

where  $l \in L_s$ . The smaller the NGMD value is, the more similar  $\mathcal{X}$  and  $\tilde{\mathcal{X}}$  are.

It is exciting how well NGMD match our intuition, see Fig. 6. The results of [18] can roughly outline the sketches but can't add textures, thus they got a high NGMD in Table 2. The results of MRF [14] and WMRF [20] are more similar to hand drawn sketches in style since these exemplar based approaches use patches of hand drawn sketches to make up the target sketch. That is why they have a smaller NGMD value. Inheriting the ability of denoising, the filtering based approach [11] is good at suppressing noises in the results. However, it is also likely for this method to over-smooth the results, which will deteriorate the texture, thus its NGMD value is higher than WMRF.

### 725 5.4. Generalization Evaluation

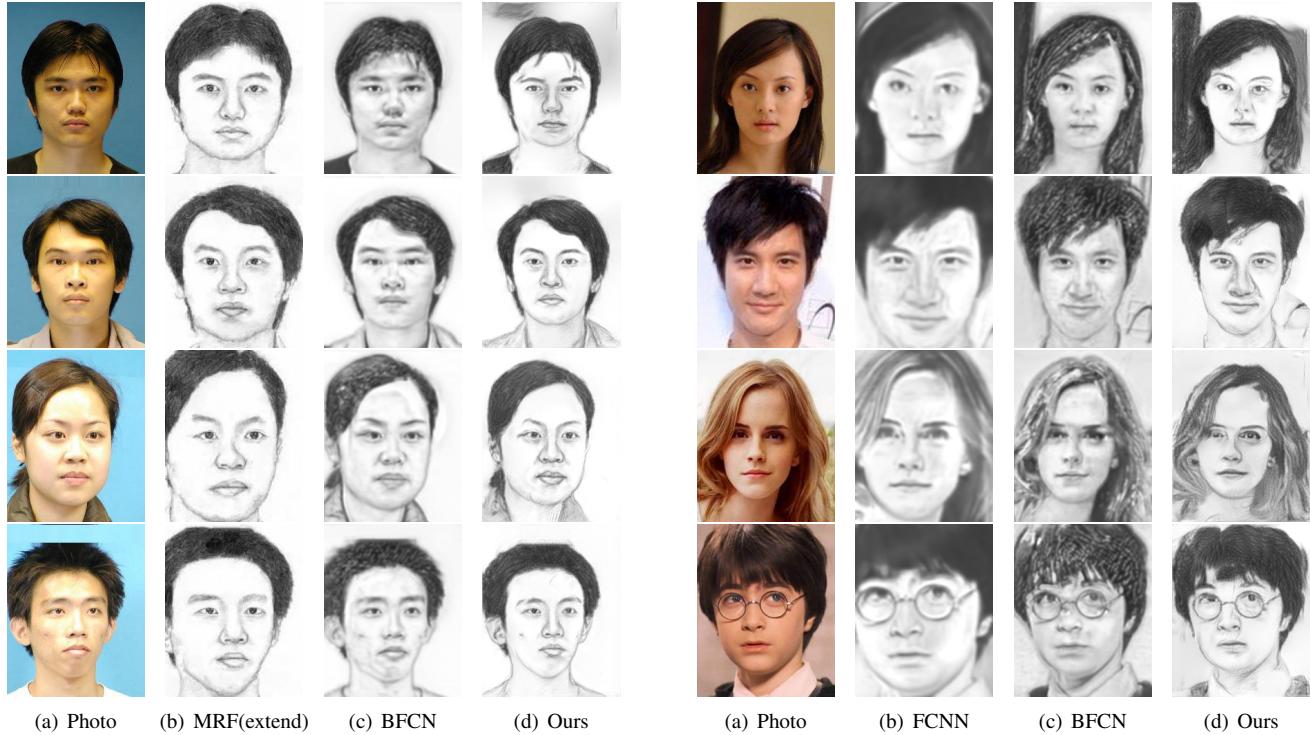
726 In this part, we will evaluate the generalization ability of 727 our model in two aspects.

**Light and Pose Invariance** As discussed in [19], light 728 and pose change may influence the result a lot. We choose 729 several photos from [19] and compare our results with 730 MRF(extended) [19] and BFCN [17]. Fig. 8 shows the 731 comparison. Our proposed method is not influenced by pose 732 and light change and can still generate good textures under lab 733 environment.

**Real World Photos** We further tested the robustness of 744 our model on some real world photos, see Fig. 9. The first 745 two rows are Chinese celebrity faces from [19], and the latter 746 two comes from the web. Since the test photo may not 747 be well aligned, we just turn off the component loss. The 748 parameters we use here are  $\alpha = 0.004$ ,  $\beta_1 = 1$ ,  $\beta_2 = 0$ . Al- 749 though the background is clutter and the positions of faces 750 are not strictly constrained. The hair style of our results are 751 still clear and sharp, while FCNN and BFCN can't produce 752 good textures.

756	Methods	AR					CUHK					810
757		conv1_1	conv2_1	conv3_1	conv4_1	conv5_1	conv1_1	conv2_1	conv3_1	conv4_1	conv5_1	811
758	FCNN	-	-	-	-	-	0.009	0.110	0.080	9.43	1.49	812
759	MRF	0.0043	0.009	0.033	0.12	0.28	0.010	0.014	0.047	0.13	0.18	813
760	WMRF	0.0053	0.027	0.085	0.19	0.29	0.010	0.052	0.052	0.27	0.19	814
761	SSD	0.0056	0.036	0.110	1.90	0.28	0.009	0.102	0.070	3.32	0.24	815
762	Ours	<b>0.0035</b>	<b>0.008</b>	<b>0.029</b>	<b>0.08</b>	<b>0.17</b>	<b>0.007</b>	<b>0.012</b>	<b>0.033</b>	<b>0.07</b>	<b>0.12</b>	816

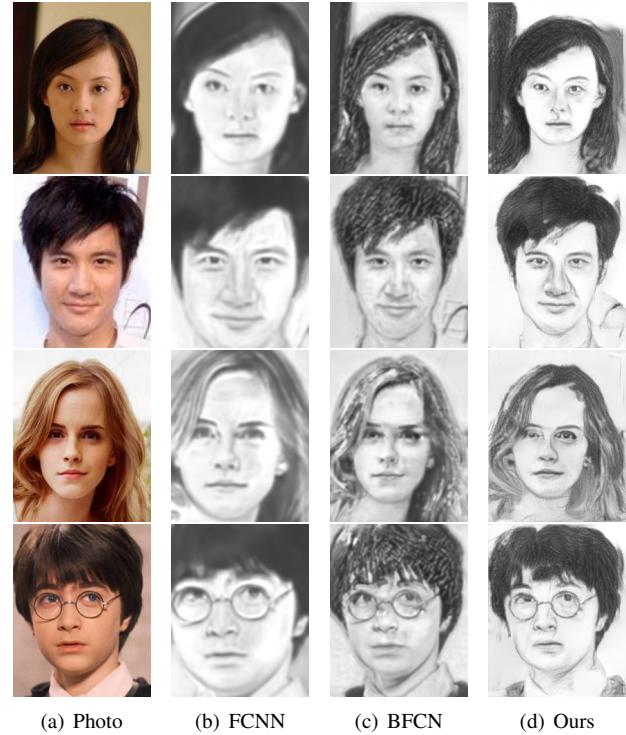
763 Table 2. Averaged NGMD value of different methods at different level on AR and CUHK datasets. A smaller NGMD indicates the  
 764 generated texture is more similar to ground truth.  
 765



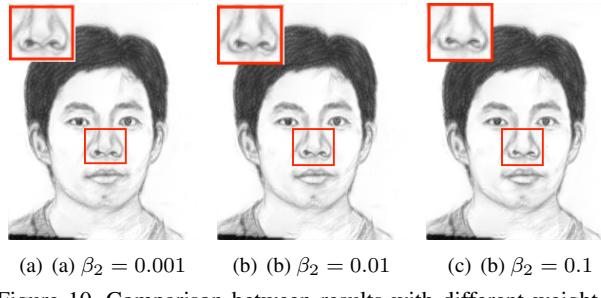
790 Figure 8. Experiment with different light and pose. Our results are  
 791 little affected by light and pose change.  
 792

## 794 5.5. Effectiveness of the model

796 The loss function we minimize during the generation  
 797 of sketches contains three terms for content, style and key  
 798 components respectively. The term  $\mathcal{L}_k$  regularizes the  
 799 results by encouraging the style extracted from the key  
 800 component regions in the training set to be placed into the key  
 801 components region of the results, which helps generate better  
 802 results around these components (see Fig. 10). To better  
 803 understand how style influences the final sketch, we  
 804 smoothly change the emphasis on style by adjusting  $\beta_1$  and  
 805  $\beta_2$  while keeping  $\alpha$  fixed. Fig. 11 indicates that the sketch  
 806 with style transferred contains more texture and is more like  
 807 a drawn sketch. The Theano implementation of the pro-  
 808 posed method takes approximately 100 seconds to generate  
 809 a sketch on a GeForce GTX TITAN X platform. The bottle



790 Figure 9. Experiments with real world photos. Our method can  
 791 still get better results under practical situations.  
 792



800 Figure 10. Comparison between results with different weight of  
 801  $\mathcal{L}_k$  regulation. With the increase of  $\beta_2$  the distortion of nose be-  
 802 comes less.  
 803

804 neck lies in the style transfer which requires feeding  $\mathcal{X}$  to  
 805 the VGG-Network to estimate targeting feature maps and to  
 806 calculate the gradient of Eq. (5), which is computationally  
 807

Figure 11. Final results for sketch synthesis with fixed  $\alpha = 0.004$  and different  $\beta_1, \beta_2$  values.

intensive.

## 6. Conclusion

This paper proposed a novel face sketch synthesis method inspired by the procedure of artists drawing sketches. In our method, the outline of the face is delineated by a content network and the style extracted from sketches drawn by artists are transferred to generate a final sketch. Quantitative evaluations on face sketch recognition and style similarity measure demonstrate the effectiveness of the proposed algorithm for face sketch synthesis and style transferring. Our future work will investigate accelerating technique to reduce the running time and achieve real time face sketch synthesis with style transfer.

## References

- [1] I. Berger, A. Shamir, M. Mahler, E. Carter, and J. Hodgins. Style and abstraction in portrait sketching. *ACM Transactions on Graphics (TOG)*, 32(4):55, 2013. [2](#)
- [2] H. Chen, Y.-Q. Xu, H.-Y. Shum, S.-C. Zhu, and N.-N. Zheng. Example-based facial sketch generation with non-parametric sampling. In *Computer Vision, 2001. ICCV 2001. Proceedings. Eighth IEEE International Conference on*, volume 2, pages 433–438. IEEE, 2001. [2](#)
- [3] T. Q. Chen and M. Schmidt. Fast patch-based style transfer of arbitrary style. *CoRR*, abs/1612.04337, 2016. [2](#)
- [4] L. Gatys, A. S. Ecker, and M. Bethge. Texture synthesis using convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 262–270, 2015. [2](#)
- [5] L. A. Gatys, A. S. Ecker, and M. Bethge. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576*, 2015. [1, 2, 3, 5](#)
- [6] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *CoRR*, abs/1502.03167, 2015. [2](#)
- [7] J. Justin, A. Alexandre, and F.-F. Li. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, pages 694–711, 2016. [2](#)
- [8] Q. Liu, X. Tang, H. Jin, H. Lu, and S. Ma. A nonlinear approach for face sketch synthesis and recognition. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 1, pages 1005–1010. IEEE, 2005. [2](#)
- [9] A. Martinez. R. benavente. the AR face database. Technical report, CVC Tech. Report, 1998. [6](#)
- [10] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. [2](#)
- [11] Y. Song, L. Bao, Q. Yang, and M.-H. Yang. Real-time exemplar-based face sketch synthesis. In *ECCV*, pages 800–813, 2014. [1, 2, 6, 7](#)
- [12] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–9, 2015. [2, 3](#)
- [13] X. Tang and X. Wang. Face sketch synthesis and recognition. In *IEEE International Conference on Computer Vision*, pages 687–694. IEEE, 2003. [2](#)
- [14] X. Wang and X. Tang. Face photo-sketch synthesis and recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(11):1955–1967, 2009. [1, 2, 3, 4, 6, 7](#)
- [15] Z. Xu, H. Chen, S.-C. Zhu, and J. Luo. A hierarchical compositional model for face representation and sketching. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(6):955–969, 2008. [2](#)
- [16] M. D. Zeiler. ADADELTA: an adaptive learning rate method. *CoRR*, abs/1212.5701, 2012. [5](#)
- [17] D. Zhang, L. Lin, T. Chen, X. Wu, W. Tan, and E. Izquierdo. Content-adaptive sketch portrait generation by decompositional representation learning. *IEEE Transactions on Image Processing*, 26(1):328–339, 2017. [1, 2, 6, 7](#)
- [18] L. Zhang, L. Lin, X. Wu, S. Ding, and L. Zhang. End-to-end photo-sketch generation via fully convolutional representation learning. In *Proceedings of the 5th ACM on International Conference on Multimedia Retrieval*, pages 627–634. ACM, 2015. [1, 2, 6, 7](#)
- [19] W. Zhang, X. Wang, and X. Tang. Lighting and pose robust face sketch synthesis. In *Computer Vision–ECCV 2010*, pages 420–433. Springer, 2010. [1, 2, 7](#)
- [20] H. Zhou, Z. Kuang, and K.-Y. K. Wong. Markov weight fields for face sketch synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1091–1097. IEEE, 2012. [1, 2, 4, 6, 7](#)