

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  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093

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.

033  
034  
035080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093

## 1. Introduction

036  
037  
038  
039  
040  
041  
042  
043  
044084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

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  
053

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

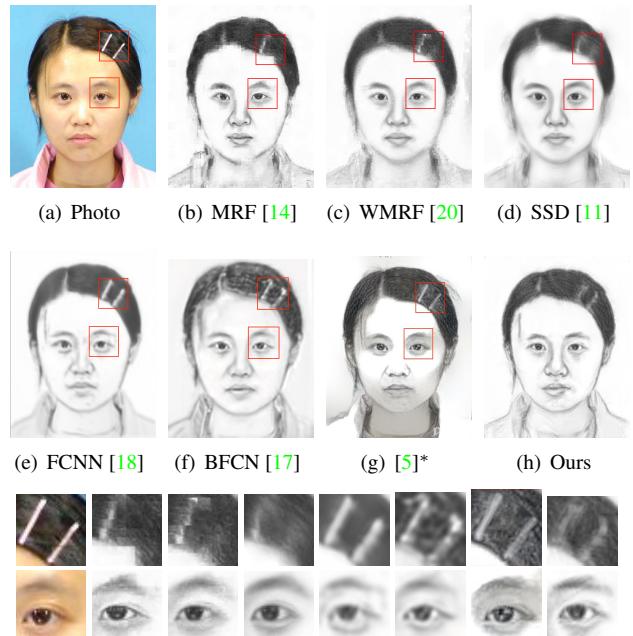


Figure 1. Face sketches generated by existing methods and the proposed method. Our method can not only preserve both hair and facial content, but also contains sharp textures. (Note that (g) is obtained from the deep art website<sup>1</sup> using the photo as content and a sketch from the training set 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(h)). 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(h)). (*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, mouth and hair. Textures and shadings  
122 are then added to regions such as hair, lips, and bridge  
123 of the nose to give the sketch a specific style. Based on  
124 the above observation, and inspired by neural style transfer  
125 [4], we propose a new framework for face sketch synthesis  
126 from a photo that overcomes the aforementioned limitations.  
127 In our method, a content image that outlines the face  
128 is generated by a feed-forward neural network, and textures  
129 and shadings are then added using a style transfer approach.  
130 Specifically, we design a new architecture of fully convolutional  
131 neural network (FCNN) composed of inception layers [12] and  
132 convolution layers with batch normalization [6] to generate  
133 the content image (see Section 4.1). To synthesize the  
134 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  
138 pyramid column features, and applied to the content image.  
139 Our approach is superior to the current state-of-the-art methods  
140 in that

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

## 147 2. Related Work

### 148 2.1. Face Sketch Synthesis

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

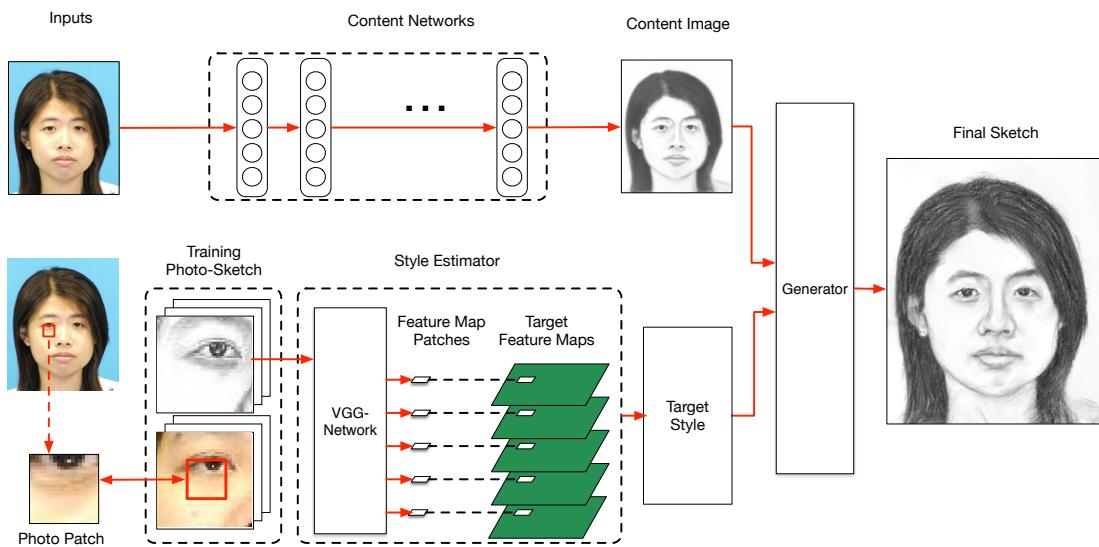
162 that the linear transformation holds better locally, and therefore  
163 they proposed a patch based method to perform sketch  
164 synthesis. In [14], a MRF based method was proposed to  
165 preserve large scale structures across sketch patches. Variants  
166 of the MRF based method were introduced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

### 179 2.2. Style Transfer with CNN

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

## 208 3. Style Representation

211 Following the work of [5], we use Gram matrices of  
212 VGG-19 [10] feature maps as our style representation. De-  
213 note the vectorized  $c$ th channel of the feature map in the  $l$ th  
214 layer of the final sketch  $\mathcal{X}$  by  $F_c^l(\mathcal{X})$ . A Gram matrix of  
215 the feature map in the  $l$ th layer is then defined by the inner

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233

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

237  
 238 products between two channels of this feature map, i.e.,  
 239

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

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

## 258 4. Methodology

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

266  
 267  
 268  
 269  
 270  
 271  
 272  
 273  
 274  
 275  
 276  
 277  
 278  
 279  
 280  
 281  
 282  
 283  
 284  
 285  
 286  
 287  
 288  
 289  
 290  
 291  
 292  
 293  
 294  
 295  
 296  
 297  
 298  
 299  
 300  
 301  
 302  
 303  
 304  
 305  
 306  
 307  
 308  
 309  
 310  
 311  
 312  
 313  
 314  
 315  
 316  
 317  
 318  
 319  
 320  
 321  
 322  
 323

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

### 336 4.1. Content Image Generation

337 The architecture of our content network is shown in  
 338 Fig. 3. Besides the test photo, we feed three extra chan-  
 339 nels containing the spatial information (i.e.,  $x$  and  $y$  coor-  
 340 dinates) and a difference of Gaussian (DoG) image into our  
 341 content network. As pointed out in [14], face sketch syn-  
 342 thesis algorithms can benefit from integrating features from  
 343 multiple resolutions. Hence, we employ an inception mod-  
 344 ule [12] for feature extraction, which concatenates features  
 345 extracted using a two-layer-inception module are then fed into  
 346 a three-layer-CNN for feature integration, where all filters  
 347 have a size of  $1 \times 1$ . Finally, the integrated features are used  
 348 to reconstruct the content image  $\mathcal{C}$  by a two-layer-CNN with  
 349 the filter size being  $3 \times 3$ . Since  $L_1$ -norm is better at pre-  
 350 serving details than  $L_2$ -norm, we use the  $L_1$ -norm between  
 351  $\mathcal{C}$  and the ground truth sketch  $S$  as the loss function in train-

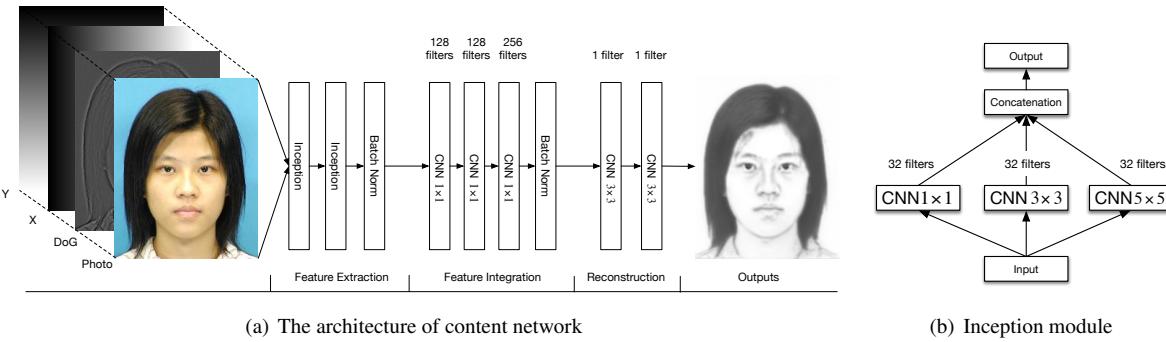


Figure 3. Illustration of the content network for generating a content image. The numbers above the building block denote the number of CNN filters. (a) The architecture of content network. (b) The inception module in (a) contains three groups of filters with different sizes.

ing our content network, i.e.,

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

where  $N$  is the number of training photos.

## 4.2. Style Estimation

As mentioned previously, there usually does not exist a single photo-sketch pair in the training set that matches all properties of the test photo. In order to estimate a target style for the final sketch  $\mathcal{X}$ , we subdivide the test photo into a grid of non-overlapping  $16 \times 16$  patches. For each test patch, similar to previous work [14, 20], we find the best matching photo patch from the photo-sketch pairs in the training set in terms of MSE (*ken: defines MSE*). A target sketch patch can then be obtained from the corresponding sketch in the photo-sketch pair containing the best matching photo patch. Instead of compositing a style image using the thus obtained target sketch patches, which may show inconsistency across neighboring patches, we adopt a feature-space approach here. We extract feature patches from the feature maps of the original sketch at 5 different layers of the VGG-Network, namely  $conv1\_1$ ,  $conv2\_1$ ,  $conv3\_1$ ,  $conv4\_1$  and  $conv5\_1$  respectively, that correspond to the target sketch patch. These feature patches have a size of  $16 \times 16$ ,  $8 \times 8$ ,  $4 \times 4$ ,  $2 \times 2$  and  $1 \times 1$  respectively (see Fig. 4). We group these five feature patches of a target sketch patch and call it a *pyramid column feature*. Finally, a target style, in the form of Gram matrices, can be computed directly from a grid of such pyramid column features.

## 4.3. Loss Function for Sketch Generation

Similar to [5], our loss function is composed of a content loss and a style loss. In addition, we introduce a component loss to enhance the key facial components. The total loss is

$$\mathcal{L}_t(\mathcal{X}) = \alpha \mathcal{L}_c + \beta_1 \mathcal{L}_s + \beta_2 \mathcal{L}_k, \quad (3)$$

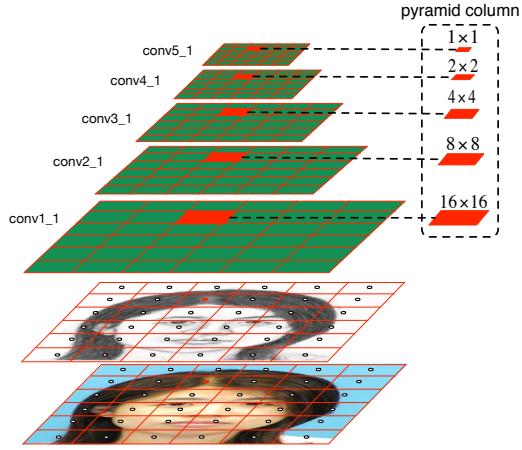


Figure 4. Illustration of the pyramid column feature. (*ken: rewrite this caption*)

where  $\alpha$ ,  $\beta$  and  $\beta_2$  are the weights for the different loss terms. We minimize the loss function by updating the target sketch  $\mathcal{X}$  without changing the VGG parameters?.

The content loss is defined by the difference between the feature map at layer  $conv1\_1$  of the synthesized sketch and that of the content image :

$$\mathcal{L}_c(\mathcal{X}) = \|F^{conv1\_1}(\mathcal{X}) - F^{conv1\_1}(\mathcal{C})\|_2^2. \quad (4)$$

The style loss is defined by the difference between the Gram matrices of the synthesized sketch and that of the target style:

$$\mathcal{L}_s(\mathcal{X}) = \sum_{l \in L_s} \frac{1}{M_l^2 N_l^2} \|G^l(\mathcal{X}) - G^l((T))\|_2^2 \quad (5)$$

where  $N_l$  denotes the number of channels of the feature map at layer  $l$ , and  $M_l$  is the product of width and height of the feature map at layer  $l$ , and  $G^l(T)$  is the Gram matrix computed from the grid of pyramid column features.

To better transfer styles of the key facial components, we employ a component loss to encourage the key component style of the final sketch to be the same as the target

432 key component style. Since all photos and sketches have  
 433 been aligned by the centers of the two eyes, the key com-  
 434 ponents lie roughly within a rectangular region  $\mathcal{R}$  with the  
 435 eyes positioned at its upper corners. Here, we define the  
 436 key component style by Gram matrices computed from fea-  
 437 ture maps corresponding to the rectangular region  $\mathcal{R}$ . The  
 438 component loss is defined as  
 439

$$\mathcal{L}_k(\mathcal{X}) = \sum_{l \in L_s} \frac{1}{\hat{M}_l^2 N_l^2} \|\hat{G}^l(\mathcal{X}) - \hat{G}^l(\mathcal{T})\|_2^2 \quad (6)$$

440 where  $\hat{G}^l$  denotes the Gram matrix computed for the re-  
 441 gular region  $\mathcal{R}$ , and  $\hat{M}_l$  is the product of width and height  
 442 of the feature map at layer  $l$  corresponding to  $\mathcal{R}$ .

#### 443 4.4. Implementation Details

444 **VGG-19 Parameters** Since the VGG-Network is origi-  
 445 nally designed for color images, while sketches are gray  
 446 scale images, we modify the first layer of VGG-Network  
 447 for gray scale images by setting the filter weights to  
 448

$$W^k = W_r^k + W_g^k + W_b^k \quad (7)$$

449 where  $W_r^k$ ,  $W_g^k$ , and  $W_b^k$  are weights of the  $k$ th filter in  
 450 the first convolutional layer for the R, G and B channels  
 451 respectively, and  $W^k$  is the weight of the  $k$ th filter in the  
 452 first convolutional layer of our modified network.  
 453

454 **Data Partition** CUHK (*ken: ref?*) has 88 training photos  
 455 and 100 test photos, and AR (*ken: ref?*) has 123 photos.  
 456 Our training set is composed of the 88 training photos of  
 457 CUHK and 100 photos from AR. When training the content  
 458 network, 10% of the training set are taken out as the validation  
 459 set. All the 188 photo-sketch pairs are used to generate  
 460 target sketch.  
 461

462 **Training the Content Network** The input photo-sketch  
 463 pairs are all resized to  $200 \times 256$  (*ken: not 288 × 288?*) and  
 464 aligned by the centers of the two eyes. A mirror padding  
 465 is carried out before the convolution operation when nec-  
 466 essary to ensure the output feature map is of the same size  
 467 as the input (*ken: ?*). Adadelta [16] is used as the optimizer  
 468 because it is stable and much faster than others.  
 469

470 **Sketch Generation** In all experiments, we resize the test  
 471 photos and the photo-sketch pairs in the training set to  $288 \times$   
 472  $288$ . The final sketch is obtained by resizing the resulting  
 473 sketch back to the original size. The size of  $\mathcal{R}$  is  $48 \times 48$ .  
 474 The weights in Eq. (3) are  $\alpha = 0.004$ ,  $\beta_1 = 1$  and  $\beta_2 = 0.1$ .  
 475 The minimization is carried out using L-BFGS. Instead of  
 476 using random noises, we use the content image as a starting  
 477 point, which will make the optimization process converge  
 478 much faster.  
 479

## 480 5. Experiments

481 We evaluate the performance of the proposed method  
 482 against other state-of-the-art methods on the CUHK student  
 483 dataset [14] and the AR dataset [9]. We compare the re-  
 484 sults of our method against other 6 methods including tra-  
 485 ditional approaches and recent deep learning models. After  
 486 discussing the disadvantages of previous quantitative eval-  
 487 uation criteria, we introduce the Normalized Gram Matrix  
 488 Difference (NGMD) as a new evaluation tool. We believe  
 489 that NGMD is more effective and will greatly promote the  
 490 sketch qualities in future work.  
 491

### 492 5.1. Style Transfer Evaluation

493 Although style transfer has shown remarkable perfor-  
 494 mance in artistic style, it can't be directly applied to face  
 495 sketch synthesis, see Fig. 1(h). The generated sketch is  
 496 greatly influenced by illumination of original photo and  
 497 doesn't look like sketch at all. To prove our assumption  
 498 in Section 3 that gram matrix captures the density distri-  
 499 buction of sketch, we replace the original RGB photo with  
 500 content image generated by our content network. We se-  
 501 lected 3 sketch styles with different amount of hairs and see  
 502 how it influences the result. Fig. 6 shows a clear rela-  
 503 tionship between the hair amount of sketch and pixel intensities  
 504 of sketch result. The facial key parts in Fig. 6(c) and 6(d)  
 505 are missed. Therefore, even with a content structure, the  
 506 original style transfer is still not suitable for elaborate task.  
 507 However, if the style image has a similar structure with test  
 508 photo, for example Fig. 6(b), the results can be quite good.  
 509 This inspires our feature patch based method. Image patch  
 510 method is not considered because it will introduce patch in-  
 511 consistency as discussed in Section 2.  
 512

### 513 5.2. Sketch Generation

514 Fig. 5 shows the comparison between our methods and  
 515 many other approaches. The first two rows are from CUHK  
 516 test set, and the last two are from AR. We can see that  
 517 our method can generate more stylistic sketches than oth-  
 518 ers. For example, in the hair part, only MRF, BFCN and  
 519 the proposed method can generate obvious textures. How-  
 520 ever, the texture of MRF is not continuous in the border part  
 521 and introduce many other artifacts, and the texture of BFCN  
 522 doesn't look like human strokes. Both WMRF and SSD in-  
 523 troduce over smoothing effect and FCNN is not able to give  
 524 clear textures. Our method can not only generate textures  
 525 for hairs and mustache but also shadings, for example the  
 526 nose part.  
 527

528 On the other hand, only BFCN and the proposed method  
 529 can handle structures decorated on the face well, for exam-  
 530 ple the glasses of the last two row. MRF, WMRF and SSD  
 531 are exemplar based method so they can't handle something  
 532 different from training set, the glass edges of them are not  
 533

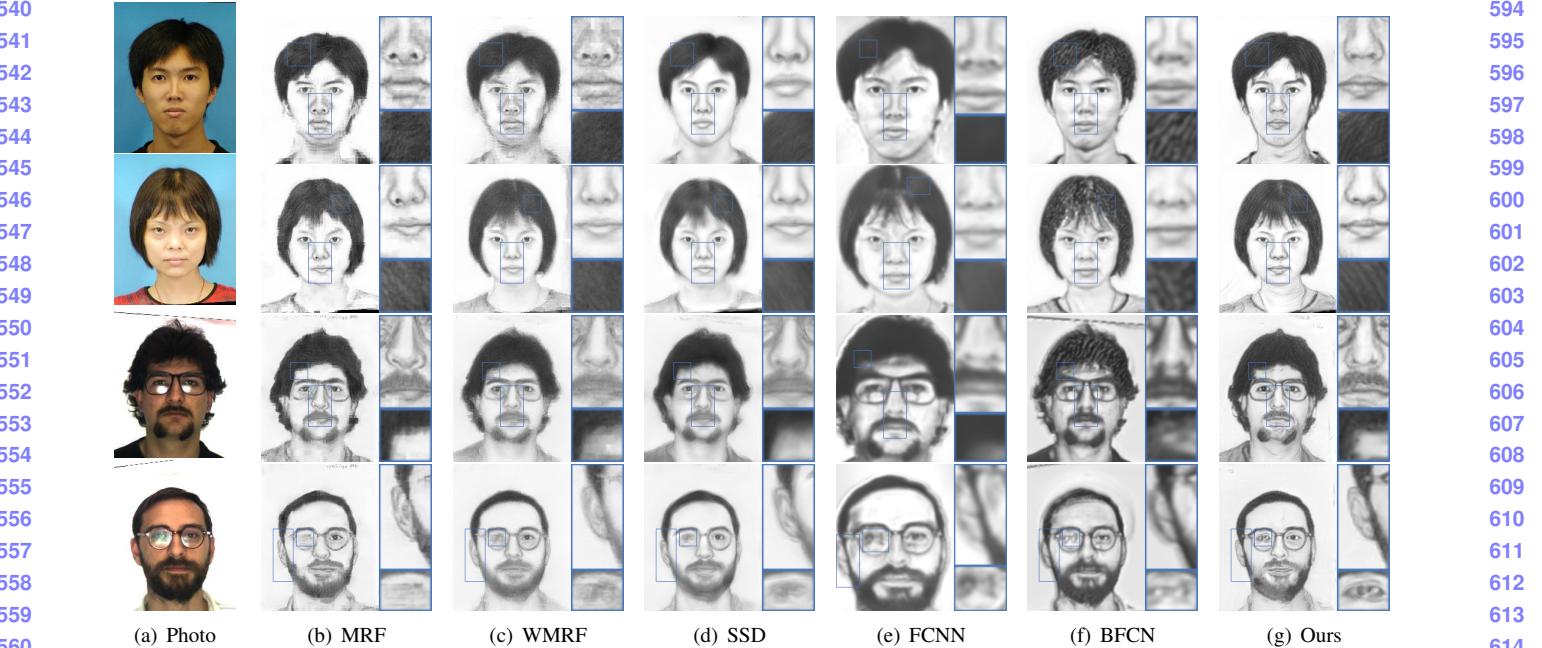


Figure 5. Examples of qualitative evaluation on CUHK (first two rows) and AR (last two rows). (a) The sketches drawn by artists. (b) MRF [14] (c) WMRF [20] (d) SSD [11] (e) FCNN [18] (f) BFCN [17] (g) Ours. The proposed method preserves more texture details, for example in the hair and nose. It is also best at keeping the origin structures of photos, such as the glasses.

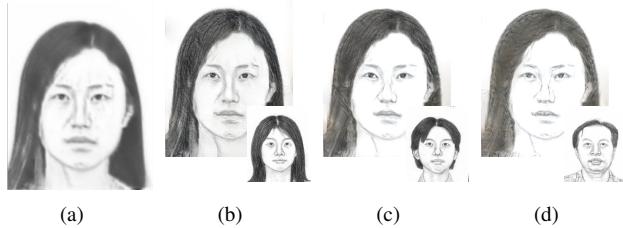


Figure 6. (a) is the content image generated by our content network. (b), (c) and (d) are generated by different styles. It can be seen that when the hair decreases the generated sketch becomes brighter.

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.

### 5.3. Quantitative Results

**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-

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%

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

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 in-

Methods	AR					CUHK				
	conv1_1	conv2_1	conv3_1	conv4_1	conv5_1	conv1_1	conv2_1	conv3_1	conv4_1	conv5_1
FCNN	-	-	-	-	-	0.009	0.110	0.080	9.43	1.49
MRF	0.0043	0.009	0.033	0.12	0.28	0.010	0.014	0.047	0.13	0.18
WMRF	0.0053	0.027	0.085	0.19	0.29	0.010	0.052	0.052	0.27	0.19
SSD	0.0056	0.036	0.110	1.90	0.28	0.009	0.102	0.070	3.32	0.24
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>

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

dicates 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  $\tilde{\mathcal{X}}$ :

$$D_s^l(\mathcal{X}, \tilde{\mathcal{X}}) = \frac{\|G^l(\mathcal{X}) - G^l(\tilde{\mathcal{X}})\|_2^2}{\|G^l(\tilde{\mathcal{X}})\|_2^2} \quad (8)$$

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. 5. 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.

#### 5.4. Generalization Evaluation

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

**Light and Pose Invariance** As discussed in [19], light and pose change may influence the result a lot. We choose several photos from [19] and compare our results with MRF(extended) [19] and BFCN [17]. Fig. 7 shows the comparison. Our proposed method is not influenced by pose and light change and can still generate good textures under lab environment.

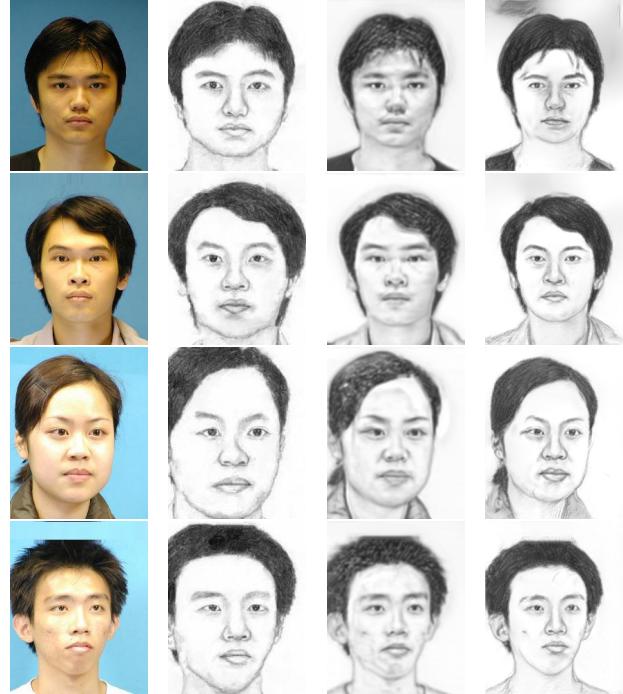
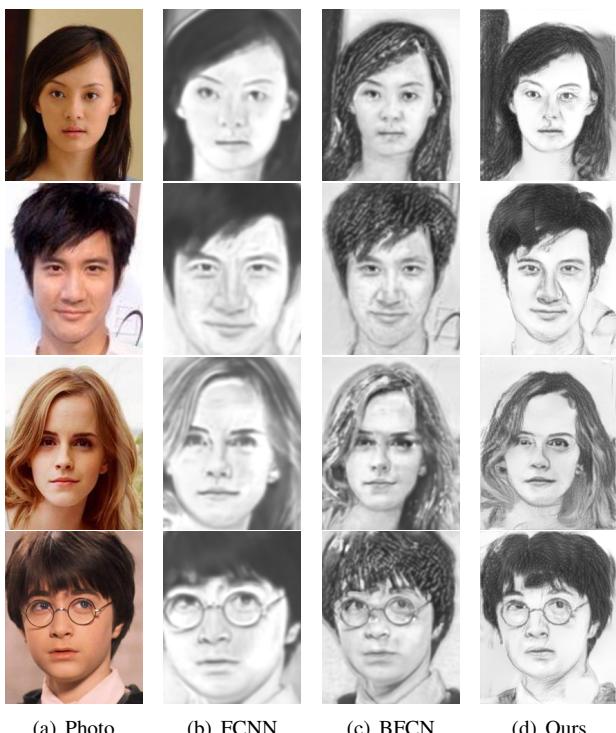


Figure 7. Experiment with different light and pose. Our results are little affected by light and pose change.

**Real World Photos** We further tested the robustness of our model on some real world photos, see Fig. 8. The first two rows are Chinese celebrity faces from [19], and the latter two comes from the web. Since the test photo may not be well aligned, we just turn off the component loss. The parameters we use here are  $\alpha = 0.004$ ,  $\beta_1 = 1$ ,  $\beta_2 = 0$ . Although the background is clutter and the positions of faces are not strictly constrained. The hair style of our results are still clear and sharp, while FCNN and BFCN can't produce good textures.

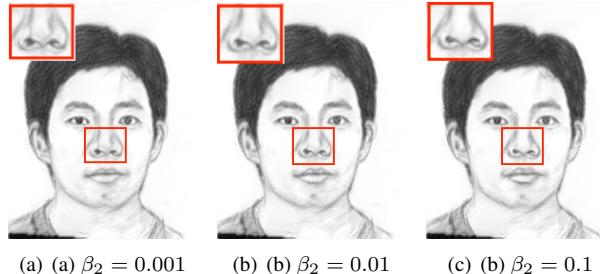
#### 5.5. Effectiveness of the model

The loss function we minimize during the generation of sketches contains three terms for content, style and key components respectively. The term  $\mathcal{L}_k$  regularizes the re-



(a) Photo      (b) FCNN      (c) BFCN      (d) Ours

Figure 8. Experiments with real world photos. Our method can still get better results under practical situations.

Figure 9. Comparison between results with different weight of  $\mathcal{L}_k$  regulation. With the increase of  $\beta_2$  the distortion of nose becomes less.

sults by encouraging the style extracted from the key component regions in the training set to be placed into the key components region of the results, which helps generate better results around these components (see Fig. 9). To better understand how style influences the final sketch, we smoothly change the emphasis on style by adjusting  $\beta_1$  and  $\beta_2$  while keeping  $\alpha$  fixed. Fig. 10 indicates that the sketch with style transferred contains more texture and is more like a drawn sketch. The Theano implementation of the proposed method takes approximately 100 seconds to generate a sketch on a GeForce GTX TITAN X platform. The bottle neck lies in the style transfer which requires feeding  $\mathcal{X}$  to the VGG-Network to estimate targeting feature maps and to

calculate the gradient of Eq. (3), which is computationally 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](#), [4](#)
- [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. [5](#)
- [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*



$\beta_1 = 0$        $\beta_1 = 10^{-3}$        $\beta_1 = 10^{-2}$        $\beta_1 = 0.1$        $\beta_1 = 1$        $\beta_1 = 10^3$   
 $\beta_2 = 0$        $\beta_2 = 10^{-4}$        $\beta_2 = 10^{-3}$        $\beta_2 = 0.01$        $\beta_2 = 0.1$        $\beta_2 = 10^2$

Figure 10. Final results for sketch synthesis with fixed  $\alpha = 0.004$  and different  $\beta_1, \beta_2$  values.

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, 5, 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

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971