

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
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
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Face Sketch Synthesis with Style Transfer using Pyramid Column Feature

Anonymous ICCV submission

Paper ID 612

Abstract

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.

1. Introduction

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.

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

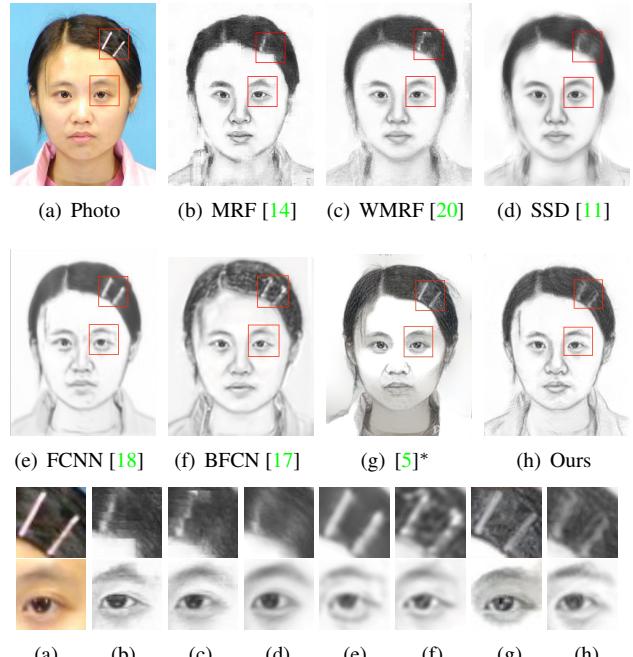


Figure 1. Face sketches generated by existing methods and the proposed method. Our method can not only preserve both hair and facial content, but also contains sharp textures. (Note that (g) is obtained from the deep art website¹ using the photo as content and a sketch from the training set as style.)

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

¹<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, 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

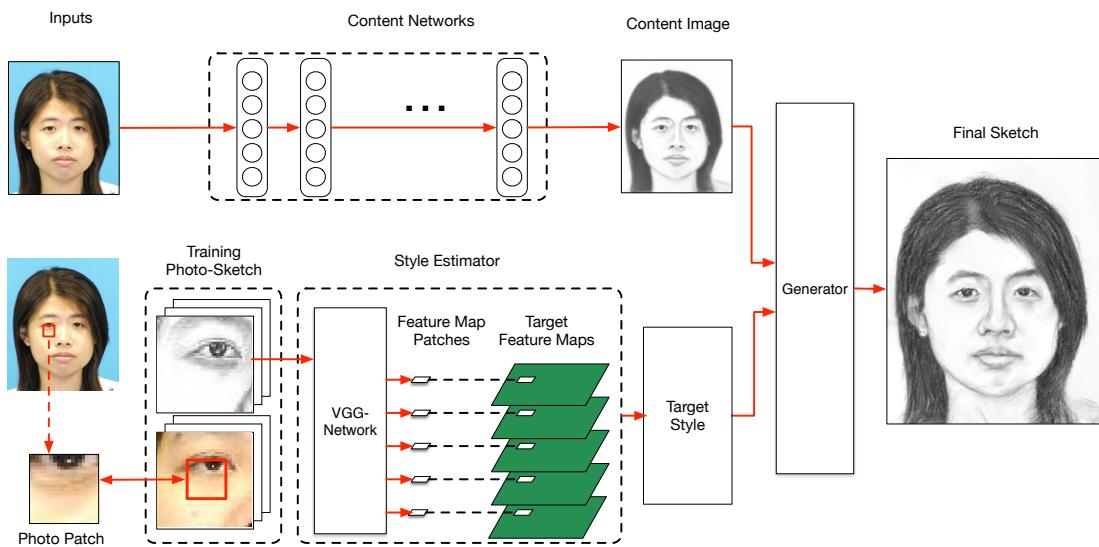
162 that the linear transformation holds better locally, and therefore
163 they proposed a patch based method to perform sketch
164 synthesis. In [14], a MRF based method was proposed to
165 preserve large scale structures across sketch patches. Variants
166 of the MRF based method were introduced in [19, 20] to
167 improve the robustness to lighting and pose, and to render
168 the ability of generating new sketch patches. In addition to
169 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. However,
202 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. Denote
213 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.

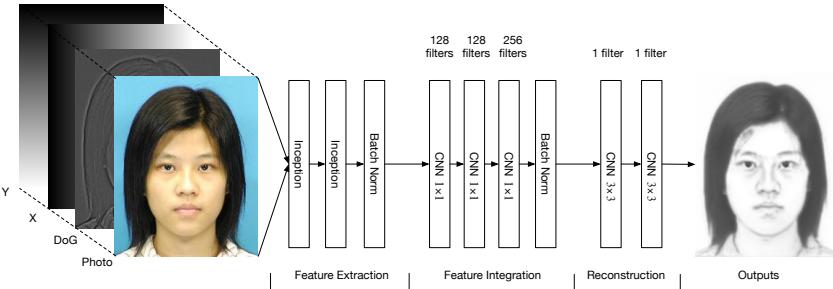
258 4. Methodology

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}

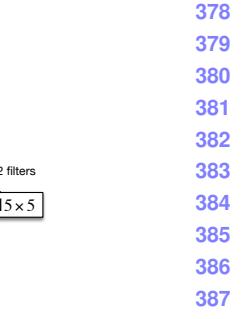
266 into a content image \mathcal{C} , which outlines the shape of the
 267 face and the key facial features such as nose, eyes, mouth
 268 and hair. The style estimator divides \mathcal{I} into a grid of non-
 269 overlapping 16×16 patches. For each test patch, it locates
 270 the most similar photo patch from the photo-sketch pairs in
 271 the training set and produce a target sketch patch from the
 272 corresponding sketch in the pair. A pyramid column feature
 273 (Section 4.2) is then computed for the target sketch patch.
 274 Finally, a target style can be computed from a grid of these
 275 pyramid column features, and a final sketch \mathcal{X} can be syn-
 276thesized by applying the target style to \mathcal{C} through neural
 277 style transfer [5].

278 4.1. Content Image Generation

279 The architecture of our content network is shown in
 280 Fig. 3. Besides the test photo, we feed three extra chan-
 281 nels containing the spatial information (i.e., x and y coor-
 282 dinates) and a difference of Gaussian (DoG) image into our
 283 content network. As pointed out in [14], face sketch syn-
 284 thesis algorithms can benefit from integrating features from
 285 multiple resolutions. Hence, we employ an inception mod-
 286ule [12] for feature extraction, which concatenates features
 287 extracted using a two-layer-inception module are then fed into
 288 a three-layer-CNN for feature integration, where all filters
 289 have a size of 1×1 . Finally, the integrated features are used
 290 to reconstruct the content image \mathcal{C} by a two-layer-CNN with
 291 the filter size being 3×3 . Since L_1 -norm is better at pre-
 292 serving details than L_2 -norm, we use the L_1 -norm between
 293 \mathcal{C} and the ground truth sketch S as the loss function in train-

324
325
326
327
328
329
330
331
332
333

(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.

334
335
336
337
338
339
340
341
342
343
344
345
346

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.

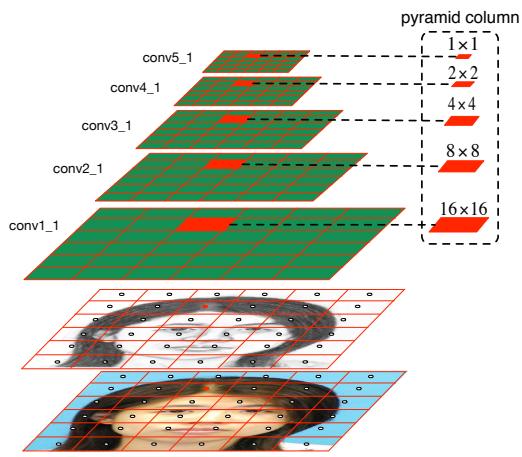
4.2. Style Estimation

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 16×16 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 **Mean Square Error (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 16×16 , 8×8 , 4×4 , 2×2 and 1×1 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*. After finding a target sketch patch, we can extract its corresponding feature patches from the original feature maps. These five patches make up the *pyramid column feature*. (*ken: rewrite this caption*)

where α , β and β_2 are the weights for the different loss terms. We minimize the loss function by updating the target sketch \mathcal{X} .

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 :

$$\mathcal{L}_c(\mathcal{X}) = \|F^{\text{conv1_1}}(\mathcal{X}) - F^{\text{conv1_1}}(\mathcal{C})\|_2^2. \quad (4)$$

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

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

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 $G^l(\mathcal{T})$ 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 key component style. Since all photos and sketches have been aligned by the centers of the two eyes, the key components lie roughly within a rectangular region \mathcal{R} with the eyes positioned at its upper corners. Here, we define the key component style by Gram matrices computed from feature maps corresponding to the rectangular region \mathcal{R} . The component loss is defined as

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

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

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

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

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 [14] (*ken: ref?*) has 88 training photos and 100 test photos, and AR [9] (*ken: ref?*) 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 the validation set. All the 188 photo-sketch pairs are used to generate target sketch.

Training the Content Network The input photo-sketch pairs are all resized to 288×288 (*ken: not 288 × 288?*) and aligned by the centers of the two eyes. A mirror padding is carried out before the convolution operation when necessary to ensure the output sketch is of the same size as the input (*ken: ?*). Adadelta [16] is used as the optimizer because it is stable and much faster than others.

Sketch Generation In all experiments, we resize the test photos and the photo-sketch pairs in the training set to 288×288 . The final sketch is obtained by resizing the resulting sketch back to the original size. The size of \mathcal{R} is 48×48 . The weights in Eq. (3) are $\alpha = 0.004$, $\beta_1 = 1$ 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 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

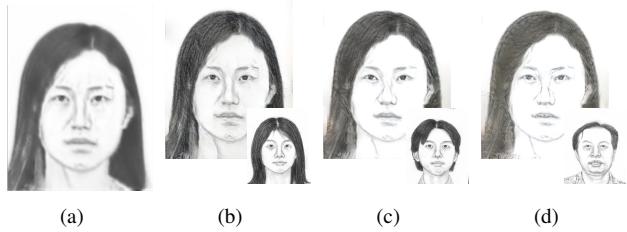


Figure 5. (a) is the content image generated by our content network. (b), (c) and (d) are generated by different styles. It can be seen that when the hair decreases the generated sketch becomes brighter. Results are generated by deepart website.

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. 5 shows a clear relationship between the hair amount of sketch and pixel intensities of sketch result. The facial key parts in Fig. 5(c) and 5(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. 5(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 inconsistency as discussed in Section 2.

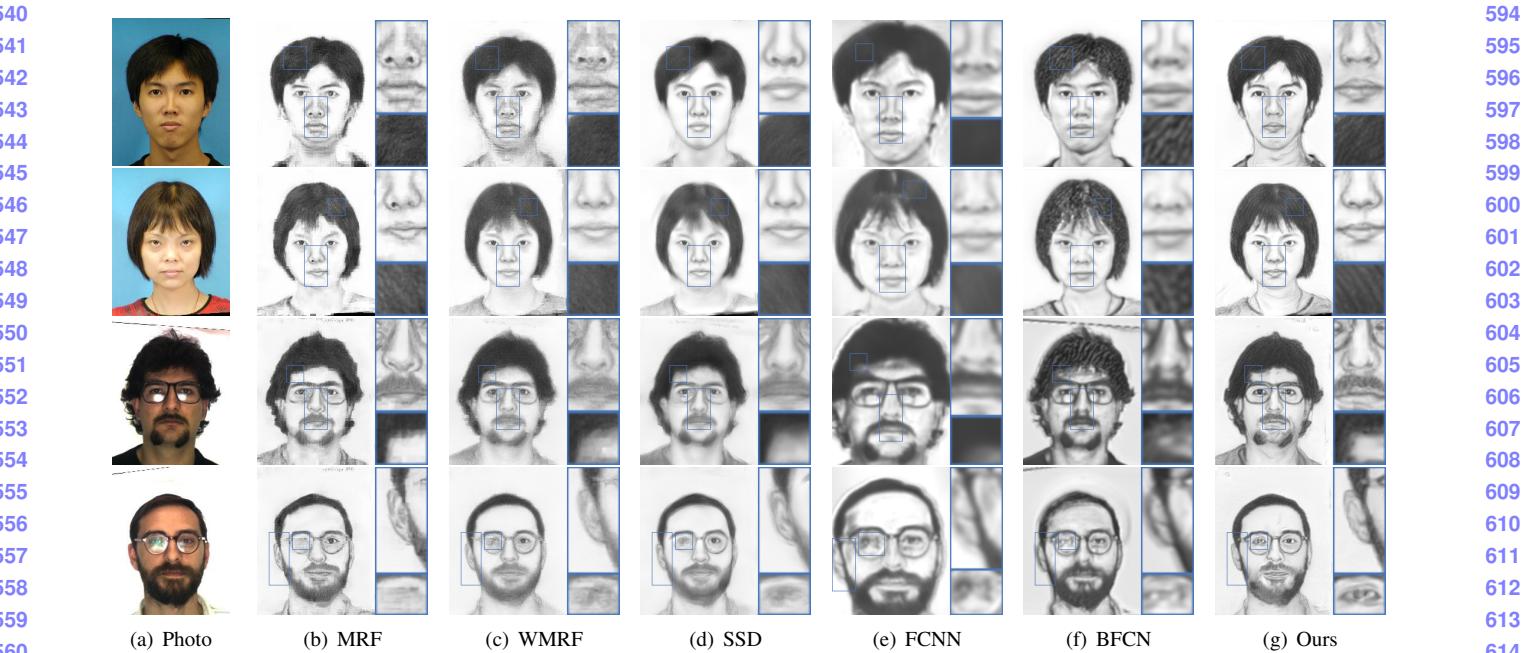


Figure 6. Examples of qualitative evaluation on CUHK (first two rows) and AR (last two rows). (a) The sketches drawn by artists. (b) MRF [14] (c) WMRF [20] (d) SSD [11] (e) FCNN [18] (f) BFCN [17] (g) Ours. The proposed method preserves more texture details, for example in the hair and nose. It is also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

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 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.

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.

5.3. 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-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 rea-

648	Methods	AR					CUHK					702
649		conv1_1	conv2_1	conv3_1	conv4_1	conv5_1	conv1_1	conv2_1	conv3_1	conv4_1	conv5_1	703
650	FCNN	-	-	-	-	-	0.009	0.110	0.080	9.43	1.49	704
651	MRF	0.0043	0.009	0.033	0.12	0.28	0.010	0.014	0.047	0.13	0.18	705
652	WMRF	0.0053	0.027	0.085	0.19	0.29	0.010	0.052	0.052	0.27	0.19	706
653	SSD	0.0056	0.036	0.110	1.90	0.28	0.009	0.102	0.070	3.32	0.24	707
654	Ours	0.0035	0.008	0.029	0.08	0.17	0.007	0.012	0.033	0.07	0.12	708

655 Table 2. Averaged NGMD value of different methods at different level on AR and CUHK datasets. A smaller NGMD indicates the
 656 generated texture is more similar to ground truth.
 657

659 sons why it is not a good criteria. First, experiment results
 660 of [18, 17] show that L_2 -norm would give blurry sketches.
 661 Second, the ground truth draw by human usually doesn't
 662 correspond exactly with photos, thus difficult for algorithms
 663 to reach the ground truth.

664 Meanwhile, human can often be able to tell whether a
 665 sketch is generated by algorithms or drawn by artists just
 666 from a quick glance without examining details. This indicates
 667 that there exists a big gap between the style of
 668 algorithm-generated sketches and that of those drawn by
 669 artists. And we need a criteria which can quantitatively
 670 measure such gap.

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

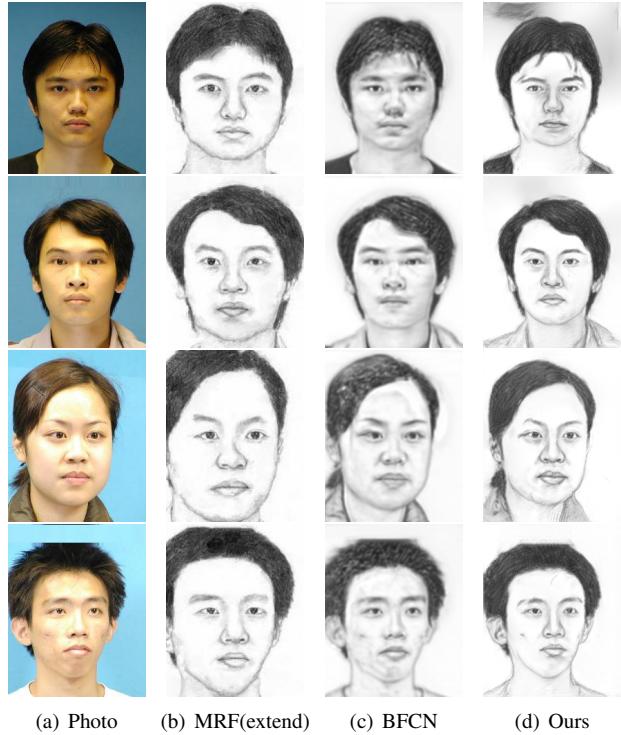
$$679 \quad D_s^l(\mathcal{X}, \tilde{\mathcal{X}}) = \frac{\|G^l(\mathcal{X}) - G^l(\tilde{\mathcal{X}})\|_2^2}{\|G^l(\tilde{\mathcal{X}})\|_2^2} \quad (8)$$

683 where $l \in L_s$. The smaller the NGMD value is, the more
 684 similar \mathcal{X} and $\tilde{\mathcal{X}}$ are.
 685

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

698 5.4. Generalization Evaluation

699 In this part, we will evaluate the generalization ability of
 700 our model in two aspects.
 701



702 Figure 7. Experiment with different light and pose. Our results are
 703 little affected by light and pose change.
 704

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

714 **Real World Photos** We further tested the robustness of
 715 our model on some real world photos, see Fig. 8. The first
 716 two rows are Chinese celebrity faces from [19], and the lat-
 717 ter two comes from the web. Since the test photo may not
 718 be well aligned, we just turn off the component loss. The
 719 parameters we use here are $\alpha = 0.004$, $\beta_1 = 1$, $\beta_2 = 0$. Al-
 720 though the background is clutter and the positions of faces
 721

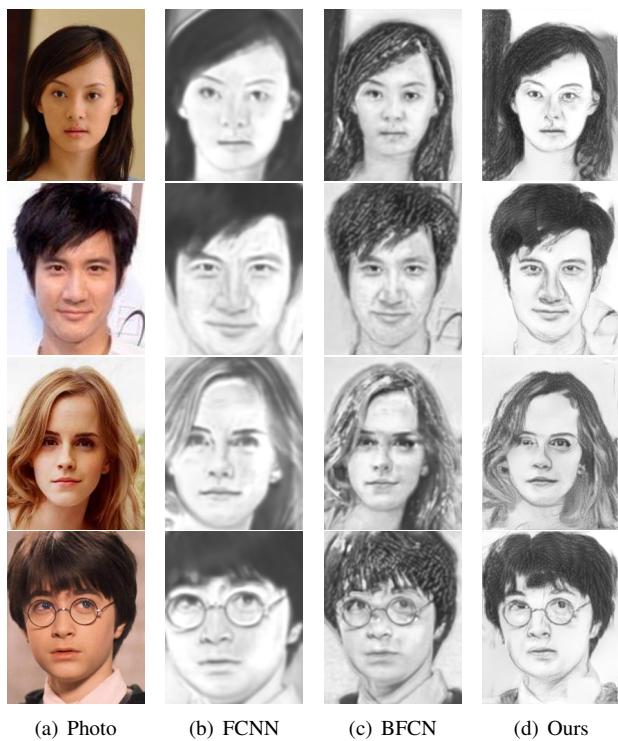


Figure 8. Experiments with real world photos. Our method can still get better results under practical situations.

are not strictly constrained. The hair style of our results are still clear and sharp, while FCNN and BFCN can't produce good textures.

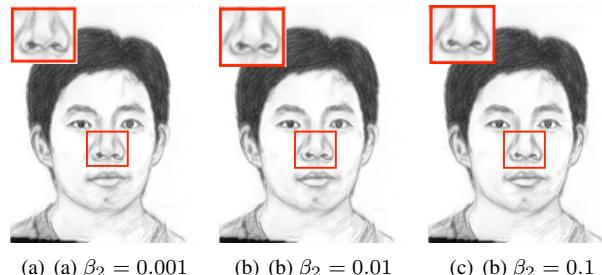


Figure 9. Comparison between results with different weight of \mathcal{L}_k regulation. With the increase of β_2 the distortion of nose becomes less.

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 results 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 β_1 and β_2 while keeping α 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](#)



$$\begin{array}{cccccc} \beta_1 = 0 & \beta_1 = 10^{-3} & \beta_1 = 10^{-2} & \beta_1 = 0.1 & \beta_1 = 1 & \beta_1 = 10^3 \\ \beta_2 = 0 & \beta_2 = 10^{-4} & \beta_2 = 10^{-3} & \beta_2 = 0.01 & \beta_2 = 0.1 & \beta_2 = 10^2 \end{array}$$

Figure 10. Final results for sketch synthesis with fixed $\alpha = 0.004$ and different β_1, β_2 values.

- [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, 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

864		918
865		919
866		920
867		921
868		922
869		923
870		924
871		925
872		926
873	$\beta_1 = 0$	927
874	$\beta_2 = 0$	928
875		929
876		930
877	[10]	931
878	K. Simonyan and A. Zisserman. Very deep convolutional	932
879	networks for large-scale image recognition. <i>arXiv preprint</i>	933
880	<i>arXiv:1409.1556</i> , 2014. 2	934
881	[11]	935
882	Y. Song, L. Bao, Q. Yang, and M.-H. Yang. Real-time	936
883	exemplar-based face sketch synthesis. In <i>ECCV</i> , pages 800–	937
884	813, 2014. 1, 2, 6, 7	938
885	[12]	939
886	C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed,	940
887	D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich.	941
888	Going deeper with convolutions. In <i>Proceedings of the IEEE</i>	942
889	<i>Conference on Computer Vision and Pattern Recognition</i> ,	943
890	pages 1–9, 2015. 2, 3	944
891	[13]	945
892	X. Tang and X. Wang. Face sketch synthesis and	946
893	recognition. In <i>IEEE International Conference on Computer Vision</i> ,	947
894	pages 687–694. IEEE, 2003. 2	948
895	[14]	949
896	X. Wang and X. Tang. Face photo-sketch synthesis and	950
897	recognition. <i>IEEE Transactions on Pattern Analysis and</i>	951
898	<i>Machine Intelligence</i> , 31(11):1955–1967, 2009. 1, 2, 3, 4, 5, 6,	952
899	7	953
900	[15]	954
901	Z. Xu, H. Chen, S.-C. Zhu, and J. Luo. A hierarchical	955
902	compositional model for face representation and sketching. <i>IEEE</i>	956
903	<i>Transactions on Pattern Analysis and Machine Intelligence</i> ,	957
904	30(6):955–969, 2008. 2	958
905	[16]	959
906	M. D. Zeiler. ADADELTA: an adaptive learning rate method.	960
907	<i>CoRR</i> , abs/1212.5701, 2012. 5	961
908	[17]	962
909	D. Zhang, L. Lin, T. Chen, X. Wu, W. Tan, and E. Izquierdo.	963
910	Content-adaptive sketch portrait generation by decompositional	964
911	representation learning. <i>IEEE Transactions on Image</i>	965
912	<i>Processing</i> , 26(1):328–339, 2017. 1, 2, 6, 7	966
913	[18]	967
914	L. Zhang, L. Lin, X. Wu, S. Ding, and L. Zhang. End-to-end	968
915	photo-sketch generation via fully convolutional representation	969
916	learning. In <i>Proceedings of the 5th ACM on International</i>	970
917	<i>Conference on Multimedia Retrieval</i> , pages 627–634.	971