# Face Sketch Synthesis Style Similarity:
# A New Structure Co-occurrence Texture Measure

Deng-Ping Fan[1], ShengChuan Zhang[2], Yu-Huan Wu[1],

Ming-Ming Cheng[1†], Bo Ren[1], Rongrong Ji[2] and Paul L Rosin[3] *

## ABSTRACT

Existing face sketch synthesis (FSS) similarity measures are sensitive to slight *image degradation* (*e.g.,* noise, blur). However, human perception of the similarity of two sketches will consider both structure and texture as essential factors and is not sensitive to slight ("pixel-level") mismatches. Consequently, the use of existing similarity measures can lead to better algorithms receiving a lower score than worse algorithms. This unreliable evaluation has significantly hindered the development of the FSS field. To solve this problem, we propose a novel and robust style similarity measure called **Scoot-measure** (Structure CO-Occurrence Texture Measure), which simultaneously evaluates "block-level" spatial structure and co-occurrence texture statistics. In addition, we further propose 4 new meta-measures and create 2 new datasets to perform a comprehensive evaluation of several widely-used FSS measures on two large databases. Experimental results demonstrate that our measure not only provides a reliable evaluation but also achieves significantly improved performance. Specifically, the study indicated a higher degree (78.8%) of correlation between our measure and human judgment than the best prior measure (58.6%). Our code will be made available.

## KEYWORDS

Face Sketch Synthesis (FSS), Structure CO-Occurrence Texture, Scoot-measure, Style Similarity

## 1 INTRODUCTION

Please take a look at Fig. 1 in which you can see several synthesized face sketches (GAN [15], FCN [37], MRF [32]). Which one do you think is closer to the ground-truth (GT) sketch drawn by the artist? While this comparison task seems trivial for humans, to date the most common measures (*e.g.,* FSIM [38], SSIM [33], VIF [21],

*[1] CCCE, Nankai University   [2] Xiamen University   [3] Cardiff University
† Ming-Ming Cheng (cmm@nankai.edu.cn) is the corresponding author.
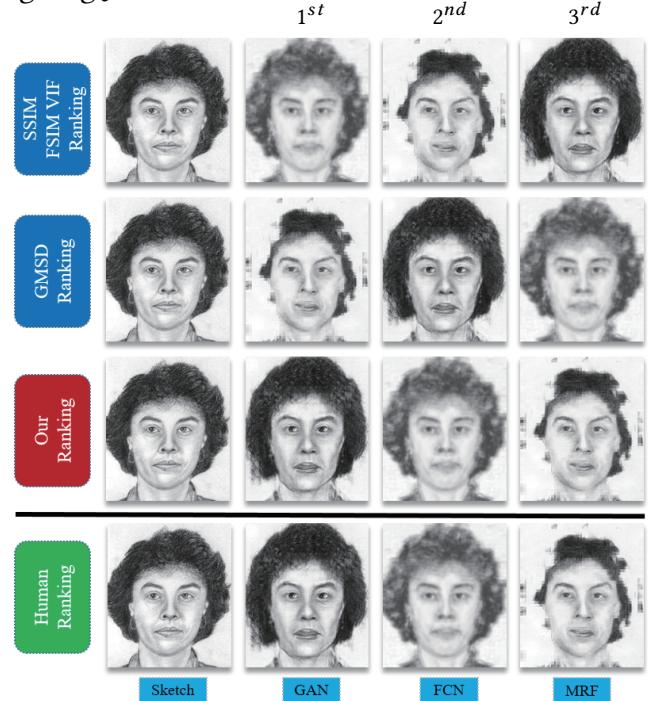Source Code: http://dpfan.net/Scoot/

Figure 1: Inaccuracy of current evaluation measures. We compare the ranking of face sketches synthesized by 3 face sketch synthesis (FSS) algorithms: GAN [15], FCN [37], and MRF [32]. According to the human ranking (last row), the GAN ranks first, followed by the FCN and MRF sketches. The GAN synthesizes the face of both structure and texture most similarly, with respect to the sketch drawn by the artist. The FCN captures the structure but lost lots of textures. The MRF almost completely destroyed the structure of the face. However, the most common measures (SSIM [33], FSIM [38], VIF [21], and GMSD [35]) fail to rank the sketches correctly. Only our measure (third row) correctly ranked the results.

and GMSD [35]) cannot reliably rank the algorithms according to their similarity to humans. This is the problem that we address in this paper, and to solve it we propose a novel measure that does much better than existing measures in terms of ranking algorithm syntheses compared to the GT sketch.

Face sketch synthesis (FSS) has been gaining popularity in the design community (*e.g.,* face artistic styles synthesis) [2, 29] and is being used for digital entertainment [11, 36, 39], and multimedia surveillance law enforcement (*e.g.,* sketch based face recognition) [4, 25]. In such applications, the comparison of a synthesized sketch

Deng-Ping Fan[1], ShengChuan Zhang[2], Yu-Huan Wu[1],
Ming-Ming Cheng[1†], Bo Ren[1], Rongrong Ji[2] and Paul L Rosin[3]

against a GT sketch is crucial in evaluating the quality of a face sketch synthesis algorithm.

To date, the most widely-used [18, 26, 29–31, 40–42, 44] evaluation measures are that of Structural SIMilarity (SSIM) [33], Visual Information Fidelity (VIF) [21], Feature SIMilarity (FSIM) [38], and Gradient Magnitude Similarity Deviation (GMSD) [35]. These measures were originally designed for *pixel-level* Image Quality Assessment (IQA) which aims to detect types of image degradation such as Gaussian blur, jpeg, and jpeg 2000 compression. Therefore these measures should be sensitive to slight (*e.g.,* pixel-level) change in image. Psychophysics [48] and prior works (*e.g.,* style for lines of drawing) [8, 12, 47] indicate that style is an essential factor in sketches. Note that human perception of the style similarity of two sketches will consider both of structure and texture [27] as the essential factor rather than being sensitive to "pixel-level" mismatches. As a consequence, current FSS measures often provide unreliable evaluation due to the different nature of their task. Consequently, a better algorithm may receive a lower score than a worse algorithm (see Fig. 1). This unreliable evaluation has significantly hindered the development of the FSS field. However, designing a reliable style similarity measure is difficult. Ideally, such a measure should:

- be **insensitive to slight mismatches** since real-world sketches drawn by artists do not exactly match each pixel to the original photos (Sec. 4.1, Sec. 4.2);

- hold **good holistic-content capture capability**, and should for example score a complete state-of-the-art (SOTA) sketch higher than the results of only preserving light strokes (Sec. 4.3);

- obtain a **high correlation** to human perception so that the sketch can be directly used in various subjective applications, especially for law enforcement and entertainment (Sec. 4.4).

As far as we know, no previous measures can meet all these goals simultaneously and provide reliable evaluation results.

To this end, we propose a new and robust style similarity measure called **Scoot-measure** (Structure CO-Occurrence Texture Measure), which simultaneously evaluates "block-level" spatial structure and co-occurrence texture statistics. We experimentally demonstrate that our measure offers a reliable evaluation and achieves significantly improved performance. Specially, the experiment indicated a higher degree (78.8%) of correlation between our measure and the human judgment than the best prior measure (58.6%). Our main contributions are:

1) We propose a simple and robust style similarity measure for face sketch synthesis that provides a unified evaluation capturing both structure and texture. We experimentally show that our measure is significantly better than most of the current widely used measures and some alternatives including GLRLM [10], Gabor [9], Canny [5], and Sobel [22].

2) To assess the quality of style similarity measures, we further propose 4 new meta-measures based on 3 hypotheses and build two new human ranked face sketches datasets. These two datasets contain 676, and 1888 synthesized face sketches, respectively. We use the two datasets to examine the ranking consistency between current measures and human judgment.

## 2 RELATED WORKS

As discussed in the introduction, a reliable similarity evaluation measure in FSS does not exist [27]. Previous researchers just used the standard measures in IQA to evaluate the similarity of FSS.

Wang *et al.* [33] proposed the **SSIM** index. Since the structural information reflects the object structure features in the scenes, SSIM computes the structure similarity, the luminance and contrast comparison using a sliding window on the local patch.

Sheikh and Bovik [21] proposed the **VIF** measure which evaluates the image quality by quantifying two kinds of information. One is obtained via the human visual system (HVS) channel, with the input ground truth and the output reference image information. The other is obtained via the distortion channel, called distortion information. And the result is the ratio of these two kinds of information.

Physiological and psychological studies found that the perceived features of human vision show great consistency with the phase consistency of Fourier series at different frequencies. Therefore, Zhang *et al.* [38] choose the phase congruency as the primary feature. Since phase congruency is not affected by the contrast, which is a key point of image quality, the image gradient magnitude is chosen as the second feature. Then they proposed a low-level feature similarity measure called **FSIM**.

Recently, Xue *et al.* [35] devised a simple measure named gradient magnitude similarity deviation (**GMSD**), where the pixel-wise gradient magnitude similarity is utilized to obtain image local quality, and the standard deviation of the overall GMS map is calculated as the final image quality index. Their measure achieves the SOTA performance compared with the other measures in IQA.

However, all the above-mentioned measures were originally designed for IQA which focuses on pixel-level image degradation, and so they are sensitive to slight resizing (see Fig. 7) of the image. Although some of them consider structure similarity, the structure is based on pixel-wise rather than pair-wise measurements. By observing sketches we found that the pair-wise co-occurrence texture is more suitable to capture the style similarity (also see our discussion in Sec. 5) than the pixel-wise texture.

## 3 SCOOT MEASURE

In this section, we introduce our novel measure to evaluate the style similarity for FSS. The important property of our measure is that it captures the pair-wise **co-occurrence texture** statistics in the "block-level" **spatial structure**. As a result, our Scoot-measure works better than the current widely-used measures. The framework is shown in Fig. 2.

### 3.1 Co-occurrence *Texture*

*3.1.1 Motivation.* To uncover the secrets of face sketch and explore quantifiable factors in style, we observed the basic principles of the sketch and noted that the "graphite pencil grades" and the "pencil's strokes" are the two fundamental elements in the sketch.

*Graphite Pencil Grades.* In the European system, "H" & "B" stand for "hard" & "soft" pencil, respectively. Fig. 4 illustrates the grade of graphite pencil from 9H-9B. Sketch images are expressed through a limited medium (graphite pencil) which provides no color. Illustrator Sylwia Bomba [24] said that "if you put your hand closer to
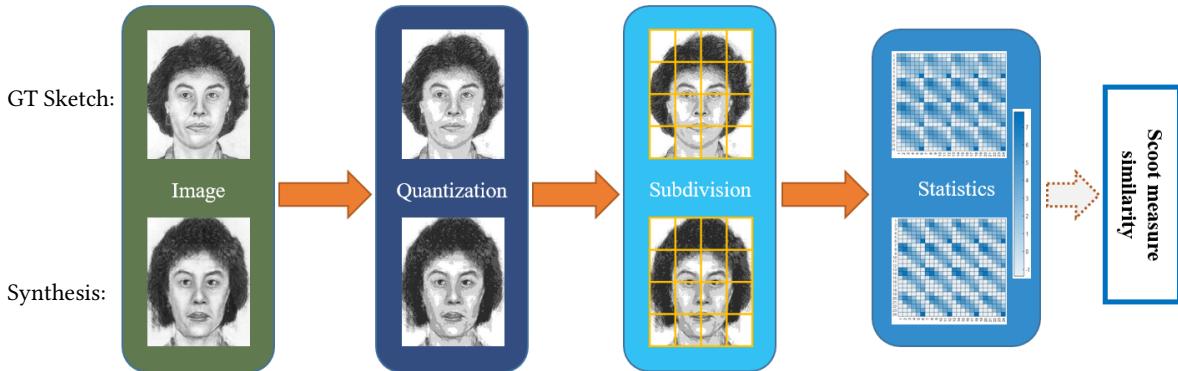
**Figure 2: The Scoot-measure framework for FSS style similarity evaluation between GT sketch and synthesis. We first quantize the two input sketches to a set of grades. Then, we divide them into blocks to consider their spatial structure. Thirdly, we extract their "block-level" co-occurrence texture statistics. Finally, the style similarity can be computed.**



(a) Outline     (b) Light     (c) Middle     (d) Dark

**Figure 3: Creating various tones of the stroke by applying different pressure on the tip. (a) is the outline of the sketch. Strokes are added to this outline by perturbing the path of the tip to generate various sketches (b, c, d). From left to right these stroke tones are light, middle (in hair regions), dark.**
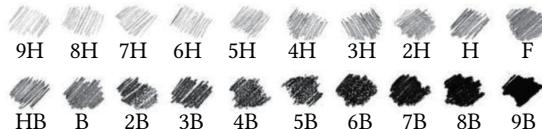


9H  8H  7H  6H  5H  4H  3H  2H  H  F

HB  B  2B  3B  4B  5B  6B  7B  8B  9B

**Figure 4: Pencil grades and their strokes.**

the end of the pencil, you have *darker* markings. Gripping further up the pencil will result in *lighter* markings." In addition, after a long period of practice, artists will form their own *fixed pressure* style. In other words, the marking of the stroke can be varied by varying the pressure on the tip (Fig. 3 (c - d)). It is important to note that different pressures on the tip will result in various types of marking which is one of the quantifiable factors called **gray tone**.

*Pencil Strokes.* Because all of the sketches are generated by moving a tip on the paper. Hence, different paths of the tip along the paper will create various stroke shapes. One example is shown in Fig. 5; different spatial distributions of the stroke have produced various textures (*e.g., sparse* or *dense*). Thus, the **stroke tone** is another quantifiable factor.
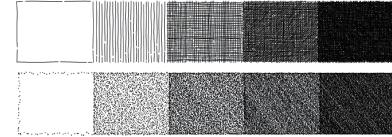


**Figure 5: Using stroke tones to indicate texture. The stroke textures used, from top to bottom, are: "cross-hatching", "stippling". The stroke attributes from left to right, are: sparse to dense. Images from [34].**

*3.1.2 Analysis.* In this section, we analyse the two quantifiable factors: gray tone and stroke tone.

**Gray tone.** To reduce the effect of slight noise and over-sensitivity to imperceptible gray tone gradient changes in sketches, intensity quantization can be introduced during the evaluation of gray tone similarity. Inspired by previous works [6], we can quantize the input sketch $I$ to $N_l$ different grades to reduce the number of intensities to be considered: $I' = \Omega(I)$.

A typical example of such quantization is shown in Fig. 6. Human beings will consistently rank (b) higher than (c) before and after quantizing the input sketches when evaluating the style similarity. Although, the quantization technology may introduces artifacts. However, our experiments (Sec. 5 of the quantized Scoot measure) show that this process can reduce sensitivity to minor intensity variations, thus achieving the best overall performance, and also reducing the computational complexity.

**Stroke tone.** Although we can efficiently capture the gray tone by quantizing the input sketch to $N_l$ grades we notice that the stroke tone (*e.g.,* dense & sparse) is also important to formulate the sketch style. Stroke tone and gray tone are not independent concepts; rather, they bear an inextricable relationship to one another. Both gray tone and stroke tone are innate properties of the sketches. The gray tone is based on the varying strokes of gray-scale in a sketch image, while the stroke tone is concerned with the spatial (statistical) distribution of gray tones.

This is demonstrated in Fig. 6, which presents the ground-truth sketch (a) and two other synthesis results (b) & (c). These two images

Deng-Ping Fan[1], ShengChuan Zhang[2], Yu-Huan Wu[1],
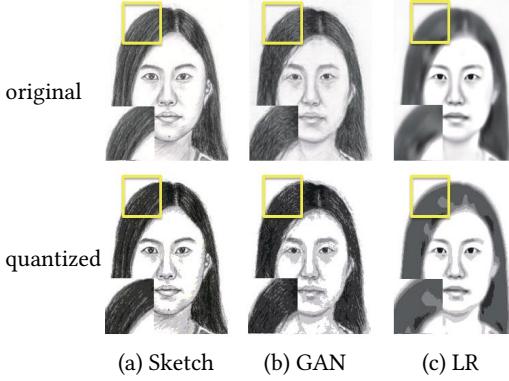Ming-Ming Cheng[1†], Bo Ren[1], Rongrong Ji[2] and Paul L Rosin[3]



**Figure 6: Different stroke attributes result in different styles. The first row contains the original sketches. The second row contains the quantized sketches. The images of the first column are the GT sketches, and the images of the second and third columns are the synthesized sketches (GAN [15], LR [31]). Texture of the hair region in (b) is more similar than that in (c) compared to the GT sketch in (a). Therefore, humans give (b) a higher score than (c).**

are the results of SOTA face synthesis algorithms. Intuitively, result (b) is better than (c), since (b) preserves the texture of the hair and details in the face while (c) presents a smooth result and has lost much of the sketch style. In other words, (b) and (c) have different stroke patterns. The co-occurrence stroke pattern of (b) has a higher similarity than the sketch in (c).

*3.1.3 Implementation.* After confirming the two quantifiable factors, we start to describe our implementation details. To simultaneously extract statistics about the "stroke tone" and their relationship to the surrounding "gray tone", we need to characterize their *spatial interrelationships*. Theoretically, the sketches are quite close to textures. As a consequence, our goal is to capture the ***co-occurrence texture*** in the sketch. Such co-occurrence properties help to capture the patterns of the face sketch (especially in the hair region). Similar to the general technology [14], we built the co-occurrence gray tones in a matrix at distance $d = (\Delta x, \Delta y)$ to encode the dependency of the gray tones. Specifically, this matrix $M$ is defined as:

$$M_{(i,j)|d} = \sum_{y=1}^{H} \sum_{x=1}^{W} \begin{cases} 1, & \text{if } I'_{x,y} = i \text{ and } I'_{(x,y)+d} = j \\ 0, & \text{otherwise} \end{cases}, \quad (1)$$

where $i$ and $j$ denote the gray value; $x$ and $y$ are the spatial positions in the given quantized sketch $I'$; $I'_{x,y}$ denotes the gray value of $I'$ at position $(x, y)$; $W$ and $H$ are the width and height of the sketch $I'$, respectively.

To extract the style features we tested the three most widely-used [13] statistics: Homogeneity, Contrast, and Energy.

**Homogeneity** reflects how much the texture changes in local regions, it will be high if the gray tone of each pixel pair is similar. The homogeneity is defined as:

$$h = \sum_{j=1}^{N_l} \sum_{i=1}^{N_l} \frac{M_{(i,j)|d}}{1 + |i - j|}, \quad (2)$$

**Contrast** represents the difference between a pixel in $I'$ and its neighbor summed over the whole sketch. This reflects that a low-contrast sketch is not characterized by low gray tones but rather by low spatial frequencies. The contrast is highly correlated with spatial frequencies. The contrast equals 0 for a constant tone sketch.

$$c = \sum_{j=1}^{N_l} \sum_{i=1}^{N_l} |i - j|^2 M_{(i,j)|d} \quad (3)$$

**Energy** measures textural uniformity. When only similar gray tones of pixels occur in an sketch ($I'$) patch, a few elements in $M$ will be close to 1, while others will be close to 0. Energy will reach the maximum if there is only one gray tone in a sketch ($I'$) patch. Thus, high energy corresponds to the sketch's gray tone distribution having either a periodic or constant form.

$$e = \sum_{j=1}^{N_l} \sum_{i=1}^{N_l} (M_{(i,j)|d})^2 \quad (4)$$

**CE (Contrast and Energy) Feature.** Our experiments (see Tab. 1) show that the combination of two statistics Contrast and Energy achieves the best average performance. In this work, we concatenate the two statistics as the CE feature to represent the style of the input sketch $I'_s$.

## 3.2 Spatial *Structure*

Note that the matrix $M$ (Eq. 1) we built is based at the "image-level". However, the weakness of this approach is that it only captures the global statistics, and the structure of the sketch is ignored. Therefore, for instance, it would be difficult to distinguish men and women with the same hairstyle.

To holistically represent the spatial structure, we follow the spatial envelope method [17] to capture the statistics from the "block-level" *spatial structure* in the sketch. Firstly, we divide the whole sketch image into a $k \times k$ grid of $k^2$ blocks[1] before extracting the CE feature. Our experiments demonstrate that the process can help to derive content information. Second, we compute the co-occurrence matrix $M$ for all blocks and normalize each matrix such that the sum of its components is 1. Finally, we concatenate the statistics of all the $k^2$ blocks into a $2k^2$ dimension vector $\overrightarrow{\Phi}(I'_s|d)$.

Note that each of the statistics is only based on a single direction (*e.g.*, $90^o$, that is d = (0, 1)), since the direction of the spatial distribution is also very important to capture the style such as "hair direction", "the direction of shadowing strokes". To exploit this observation for efficiently extracting the *stroke direction* style, we compute the average feature $\overrightarrow{\Psi}(I'_s)$ of four orientations[2] vectors to capture more directional information:

$$\overrightarrow{\Psi}(I'_s) = \frac{1}{4} \sum_{i=1}^{4} \overrightarrow{\Phi}(I'_s|d_i), \quad (5)$$

where $d_i$ denotes the $i$th direction and the average CE feature $\overrightarrow{\Psi}(I'_s) \in \mathbb{R}^{2k^2}$.

---

[1] We set $k = 4$ to achieve the best performance in our all experiments.

[2] Due to the symmetry of the co-occurrence matrix $M(i, j)$, the statistical features in 4 orientations are actually equivalent to 8 neighbors direction at 1 distance. Empirically, we set 4 orientations $d_i \in \{(0, 1), (-1, 1), (-1, 0), (-1, -1)\}$ to achieve the best performance.

### 3.3 Scoot measure

Having obtained the style feature vectors of the GT sketch $Y$ and synthesis sketch $X$, we use the Euclidean distance between their feature vectors to evaluate their style similarity. Thus, our style **Scoot-measure**[3] can be defined as:

$$E_s = \frac{1}{1 + \left\| \overrightarrow{\Psi}(X_s') - \overrightarrow{\Psi}(Y_s') \right\|_2},$$ (6)

where $\|\cdot\|_2$ denotes $l_2$-norms. $X_s'$, $Y_s'$ denote the quantized $X_s$, $Y_s$, respectively. $E_s = 1$ corresponds to identical style. Using this measure to evaluate the three synthesized face sketches in Fig. 1, we can correctly rank the sketches consistent with the human ranking.

## 4 EXPERIMENT

In this section, we compare our Scoot-measure with 8 measures including 4 popular measures (FSIM, SSIM, VIF, GMSD) and 4 alternative baseline measures (GLRLM [10], Gabor [9], Canny [5], Sobel [22]) which related with texture or edge. For the four baseline measures, we only replace our CE feature (describe in Sec. 3.1.3) with these features.

**Meta-measure.** To test the quality of our measure, we adopt the *meta-measure* methodology. The meta-measures consist of assuming some plausible hypotheses about the results and assessing how well each measure reflects these hypotheses [7, 20]. We design 4 novel meta-measures based on our plausible hypotheses (Sec. 1) and the experiment results are summarized in Tab. 2.

**Dataset.** All of the meta-measures were performed on 2 widely used databases: **CUFS** [32] and **CUFSF** [43]. The CUFS database includes 606 face photo-sketch pairs consisting of 3 sub-sections: CUHK Student (188 pairs), XM2VTS (295 pairs) and Purdue AR (123 pairs). The CUFSF database contains 1194 pairs of persons from the FERET database [19]. There is one face photo and one GT sketch drawn by the artist for each person in both CUFS and CUFSF databases. Following [26], we use 338 pairs (CUFS) and 944 pairs (CUFSF) test set images to conduct our experiments.

**Sketch Synthesis Results.** Sketch synthesis results were generated for each test image using 10 SOTA models (FCN [37], GAN [15], LLE [16], LR [31], MRF [32], MWF [45], RSLCR [26], SSD [23], DGFL [46], BFSS [28]).

### 4.1 Meta-measure 1: Stability to Slight Resizing

The first meta-measure specifies that the rankings of synthetic results should not change much with slight changes in the GT sketch. Therefore, we perform a slight 5 pixels downsizing of the GT by using nearest-neighbor interpolation.

Fig. 7 gives an example. The hair of the GT sketch in (b) drawn by the artist has a slight size discrepancy compared to the photo (a). We observe that about 5 pixels deviation (Fig. 7 c) in the boundary is common. While the two sketches (e) & (f) are almost identical, commonly-used measures including SSIM, VIF, and GMSD switched the ranking of the two synthesized results (g, h) when using (e) or (f). However, our measure consistently ranked (g) higher than (h).

Here, we applied the $\theta = 1 - \rho$ [3] measure to test the measure ranking stability before and after the GT downsizing was performed.

(a) Photo    (b) GT    (c) Differenced    (d) downsized
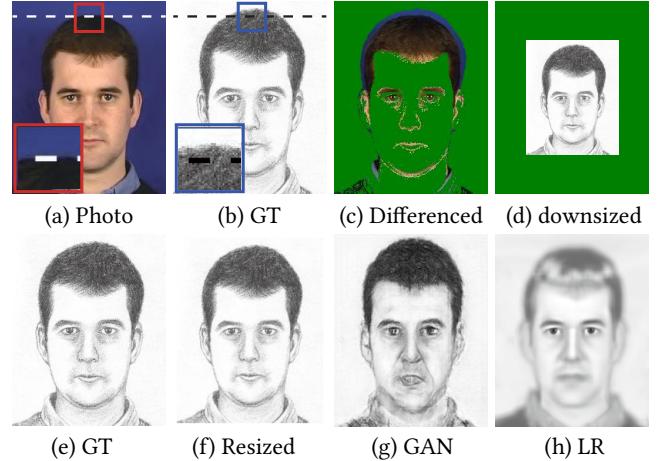
(e) GT    (f) Resized    (g) GAN    (h) LR

**Figure 7: Meta-measure 1. (a) is the original photo. (b) is the GT sketch drawn by the artist. (c) is the boundary difference between (a) & (b). (d) is the downsized image from (b). The ranking of an evaluation measure should be insensitive to slight resizing of the GT. While GT (e) & 5 pixels scaled-down GT (f) differ slightly, widely-used measures (*e.g.,* VIF, SSIM, and GMSD) switched the ranking order of the two synthesized sketches (g) & (h) when using the different versions of GT referenced (e, f). In contrast, our Scoot measure consistently ranked (g) higher than (h).**
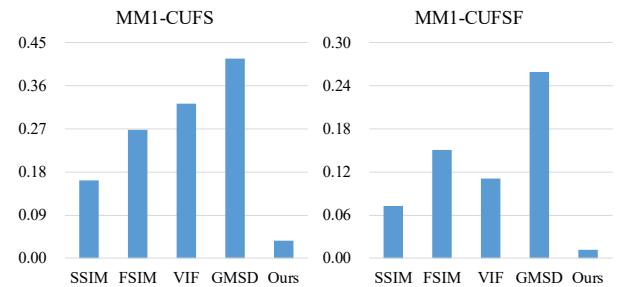


**Figure 8: Meta-measure 1 results: The lower the result is, the more stable an evaluation measure is to slight resizing.**

Fig. 8 and Tab. 2 show the results: the lower the result is, the more stable an evaluation measure is to slightly downsizing. We can see a significant ($\approx$ 88% and 92%) improvement over the existing SSIM, FSIM, GMSD, and VIF measures in both the CUFS and CUFSF databases. These improvements are mainly because our evaluation measure considers "block-level" features rather than "pixel-level".

### 4.2 Meta-measure 2: Rotation Sensitivity

In real-wold situations, sketches drawn by artists may also have slight rotations compared to the original photographs. Thus, our second meta-measure verifies the sensitivity of GT rotation for the evaluation measure. We did a slight counter-clockwise rotation ($5^o$) for each GT sketch. Fig. 9 shows an example. When the GT (a) is switched to the slightly rotated GT (b), the ranking results should not change much.

Deng-Ping Fan[1], ShengChuan Zhang[2], Yu-Huan Wu[1],
Ming-Ming Cheng[1†], Bo Ren[1], Rongrong Ji[2] and Paul L Rosin[3]



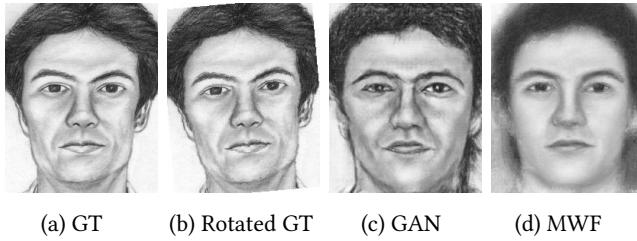(a) GT     (b) Rotated GT     (c) GAN     (d) MWF

**Figure 9: Meta-measure 2. The ranking of an evaluation measure should be insensitive to slight rotation of the GT. While GT (a) and $5^o$ rotated GT (b) has the same style, all of the current measures switched the ranking order of the two synthesized sketches (c) & (d), relying on the different GT referenced. Oppositely, our Scoot measure consistently ranked (c) higher than (d).**
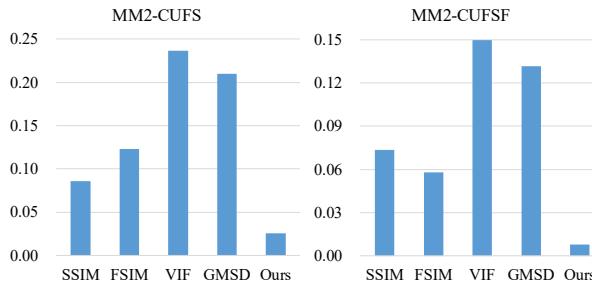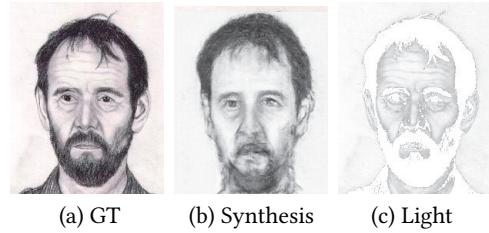


(a) GT     (b) Synthesis     (c) Light

**Figure 11: Meta-measure 3. We use a simple threshold to separate the GT sketch (a) into dark strokes and light strokes (c). From the style similarity perspective, a good evaluation measure should prefer the SOTA synthesized sketch (b) which contains the main texture and structure over the result only contains light strokes and has lost the main facial features (c).**



**Figure 10: Meta-measure 2 results: The lower the result is, the more stable an evaluation measure is to slight rotation.**
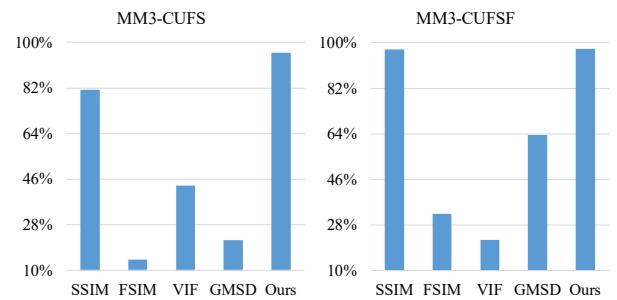


**Figure 12: Meta-measure 3 results: The higher the result, the stronger is the holistic content capture capability. Our measure achieves the best performance in both the CUFS and CUFSF databases.**

We got the ranking results for each measure by using GT sketches and slightly rotated GT sketches separately and applied the same measure ($\theta$) as meta-measure 1 mentioned to evaluate the rotation sensitivity. The sensitivity results are shown in Fig. 10, Tab. 2. Thanks to our use of "block-level" statistics, our measure again significantly outperforms the current measures over the CUFS and CUFSF databases.

### 4.3 Meta-measure 3: Content Capture Capability

The third meta-measure describes that a good evaluation measure should not only evaluate the style similarity but also can capture the content similarity. Fig. 11 presents an example. We expect that an evaluation measure should prefer the SOTA synthesized result over the **light strokes**[4] results. To our surprise, only our measure gives the correct order.

We use the *mean score* of 10 SOTA synthesis algorithms (FCN, GAN, LLE, LR, MRF, MWF, RSLCR, SSD, DGFL, BFSS). The mean score is robust to situations in which a certain model generates a poor result. We recorded the number of times the mean score of SOTA synthesized algorithms is higher than a light stroke's score.

---

[4] To test the third meta-measure, we use a simple threshold of gray scale (*e.g.,* 170) to separate the sketch (Fig. 11 GT) into darker strokes & lighter strokes. The lighter strokes image loses the main texture features of the face (*e.g.,* hair, eye, beard), resulting in an incomplete sketch.

The results are shown in Fig. 12 & Tab. 2. In the CUFS database we can see a great improvement over the most widely-used measures. A slight improvement is also achieved for the CUFSF database.

The reason why other measures fail in this meta-measure is shown in Fig. 11. In terms of pixel-level matching, it is obvious that the regions where dark strokes are removed are totally different from the corresponding parts in (a). But at other positions the pixels are identical to the GT. Previous measures only consider "pixel-level" matching, and will rank the light strokes sketch higher. However, the synthesized sketch (b) is better than the light one (c) in terms of both style and content.

### 4.4 Meta-measure 4: Human Judgment

The fourth meta-measure specifies that the ranking result according to an evaluation measure should agree with the human ranking. As far as we know, there is no such human ranking sketch database publicly available.

**Source images.** As mentioned in Sec. 4, our source sketch synthesized results are from the two widely-used datasets: CUFS and CUFSF and 10 SOTA sketch synthesis algorithms. Thus, we have 3380 (10 × 338) and 9440 (10 × 944) synthesized results for CUFS and CUFSF, respectively.

**Figure 13: Meta-measure 4. Sample images from our human ranked database. The first row is the GT sketch, followed by the first and second ranked synthesis result. Please refer the supplementary material for more images.**

**Human judgment.** For each photo of the CUFS and CUFSF databases, we have 10 sketches synthesized by 10 SOTA algorithms. Due to the similarity of the synthesized results of some algorithms, it is difficult for viewers to judge the ranking. Following [27], we trained the viewers to assess synthesized sketch quality based on two criterion: texture similarity and content similarity. To minimize the ambiguity of human ranking, we asked 5 viewers to select 2 synthesized sketches for each photo through the following 3 stage processes:

**i)** We let the first group of viewers select 4 out of 10 sketches for each photo. The 4 sketches should consist of two good and two bad ones. Thus we are left with 1352 ($4 \times 338$) and 3776 ($4 \times 944$) sketches for CUFS and CUFSF, respectively.

**ii)** For the 4 sketches in each photo, the second group of viewers is further asked to select 3 sketches for which they can rank them easily. Based on the voting results of viewers, we pick out the 3 most selected sketches.

**iii)** The last group of viewers are asked to pick out a pair of sketches that are most obvious to rank. Note that we have 5 volunteers involved in the whole process for cross-checking the ranking. Finally, we create two new human ranked datasets: **CUFS-Ranked** and **CUFSF-Ranked** (see the *supplementary* for all images in these two datasets) [5]. Fig. 13 shows the human ranked examples.

Here, we examined the the consistency between human ranking and evaluation measure ranking. Results are given in Fig. 14, which shows a significant (about 26.3%) improvement over the best prior measure in CUFS. This improvement is due to our consideration of style similarity which human perception considers as an essential factor when evaluating sketches.

## 5 DISCUSSION

In this section, we will discuss the combination of statistics in Sec. 5.1 and other alternative features in Sec. 5.2. These discussions

---

[5]They include 776 & 1888 human ranked images, respectively. The two datasets will be made publicly available.
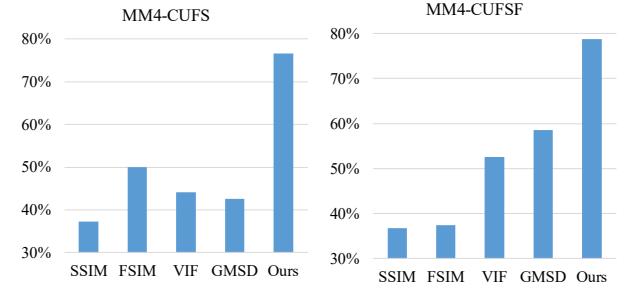


**Figure 14: Meta-measure 4 results: The higher the result, the more consistency an evaluation measure is to human judgment. Our measure achieves the best performance in both CUFS and CUFSF databases.**

**Table 1: Quantitative comparison of different combination of 3 statistical features on our 4 meta-measures. MM: meta-measure. The best balanced performance is highlighted in bold. These differences are all statistically significant at the $\alpha < 0.05$ level.**

| Measure | CUFS | | | | CUFSF | | | |
|---|---|---|---|---|---|---|---|---|
| | MM1 | MM2 | MM3 | MM4 | MM1 | MM2 | MM3 | MM4 |
| $h$ | 0.007 | 0.005 | 61.5% | 77.5% | 0.003 | 0.003 | 79.1% | 77.8% |
| $e$ | 0.200 | 0.104 | 98.5% | 73.1% | 0.044 | 0.026 | 99.2% | 77.4% |
| $c$ | 0.010 | 0.007 | 54.4% | 74.6% | 0.009 | 0.006 | 64.7% | 73.4% |
| $c + h$ | 0.011 | 0.007 | 60.1% | 74.6% | 0.007 | 0.005 | 78.1% | 73.7% |
| **$c + e$** | **0.037** | **0.025** | **95.9%** | **76.3%** | **0.012** | **0.008** | **97.5%** | **78.8%** |
| $h + e$ | 0.156 | 0.088 | 97.9% | 75.7% | 0.030 | 0.017 | 98.8% | 80.3% |
| $h + e + c$ | 0.034 | 0.024 | 95.9% | 76.3% | 0.011 | 0.008 | 97.4% | 78.7% |

indicate that the pair-wise "co-occurrence" properties are quite important for the style similarity evaluation for FSS.

### 5.1 Statistical Feature Selection

According to Sec. 3, we considered three widely-used features of co-occurrence matrices: homogeneity ($h$), contrast ($c$) and energy ($e$). In order to achieve the best performance, we need to explore the best combination of these statistical features. We have applied our four meta-measures to test the performance of our measure using each single feature, each feature pair and the combination of all three features.

The results are shown in Tab. 1. All possibilities ($h, e, c, c + e, h + e, c + h, h + e + c$) perform well in MM4 (human judgment). $h$ and $c$ are insensitive to resizing (MM1) and rotation (MM2), while they are not good at content capture (MM3). $e$ is the opposite compared to $h$ and $c$. Thus, using a single feature is not good. The results of combining two features show that if $h$ is combined with $e$, the sensitivity to resizing and rotating will still be high, while partially overcoming the weakness of $e$. The performance of $h + e + c$ shows no improvement compared to the combination of $c + e$. Previous work in [1] also found energy and contrast to be the most efficient features for discriminating textural patterns. Therefore, we choose contrast and energy as our final combination "CE" feature to extract the style features.

Deng-Ping Fan[1], ShengChuan Zhang[2], Yu-Huan Wu[1],
Ming-Ming Cheng[1†], Bo Ren[1], Rongrong Ji[2] and Paul L Rosin[3]

**Table 2: Quantitative comparison of different evaluation measures on our 4 meta-measures. MM: meta-measure. The top measure for each MM are shown in boldface. These differences are all statistically significant at the $\alpha < 0.05$ level.**

| Measure | CUFS (3380 synthesized images) | | | | CUFSF (9440 synthesized images) | | | |
|---|---|---|---|---|---|---|---|---|
| | MM1 | MM2 | MM3 | MM4 | MM1 | MM2 | MM3 | MM4 |
| SSIM [33] | 0.162 | 0.086 | 81.4% | 37.3% | 0.073 | 0.074 | 97.4% | 36.8% |
| FSIM [38] | 0.268 | 0.123 | 14.2% | 50.0% | 0.151 | 0.058 | 32.4% | 37.5% |
| VIF [21] | 0.322 | 0.236 | 43.5% | 44.1% | 0.111 | 0.150 | 22.2% | 52.8% |
| GMSD [35] | 0.417 | 0.210 | 21.9% | 42.6% | 0.259 | 0.132 | 63.6% | 58.6% |
| **Scoot(CE)** | **0.037** | **0.025** | **95.9%** | **76.3%** | **0.012** | **0.008** | **97.5%** | **78.8%** |

**Table 3: Quantitative comparison of different evaluation measures on our 4 meta-measures. MM: meta-measure. Scoot(CE)[†] indicates the non-quantized setting. The top measure for each MM are shown in boldface. These differences are all statistically significant at the $\alpha < 0.05$ level.**

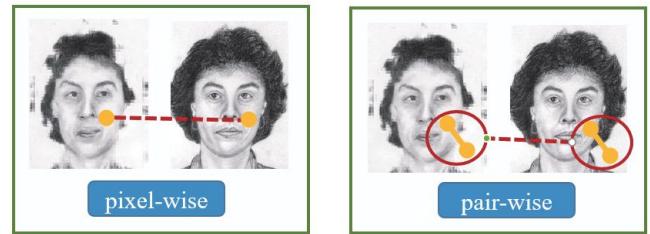| Measure | CUFS | | | | CUFSF | | | |
|---|---|---|---|---|---|---|---|---|
| | MM1 | MM2 | MM3 | MM4 | MM1 | MM2 | MM3 | MM4 |
| Scoot(Canny) | 0.086 | 0.078 | 33.7% | 27.8% | 0.138 | 0.146 | 0.0% | 0.1% |
| Scoot(Sobel) | 0.040 | 0.037 | 0.0% | 32.8% | 0.048 | 0.044 | 0.0% | 52.6% |
| Scoot(GLRLM) | 0.111 | 0.111 | 18.6% | 73.7% | 0.125 | 0.079 | 64.6% | 68.0% |
| Scoot(Gabor) | 0.062 | 0.055 | 0.0% | 72.2% | 0.089 | 0.043 | 19.3% | **80.9%** |
| Scoot(CE)[†] | **0.022** | **0.014** | 48.8% | 73.1% | 0.012 | **0.007** | 92.5% | 79.4% |
| **Scoot(CE)** | 0.037 | 0.025 | **95.9%** | **76.3%** | **0.012** | 0.008 | **97.5%** | 78.8% |

## 5.2 Alternative Features

As described in Sec. 3, sketches are quite close to textures. There are many other texture & edge-based features (*e.g.,* GLRLM [10], Gabor [9], Canny [5], Sobel [22]). Here, we select the most wide-used features as candidate alternatives to replace our "CE" feature. For GLRLM, we select all five statistics mentioned in the original version. Results are shown in Tab. 3. Gabor and GLRLM are texture features, while the other two are edge-based. With all the texture features our measure provides a high consistency with human ranking (MM4). However, GLRLM does perform well according to MM1 & 2 & 3. Gabor is reasonable in terms of MM1 & 2, but fails in MM3. For edge-based features, Canny fails according to all meta-measures. Sobel is very stable to slight resizing (MM1) or rotating (MM2), but fails to capture content (MM3) and is not consistent with human judgment (MM4).

**New insight.** Note that the alternative features *do not consider the pair-wise "co-occurrence" texture properties* (see Fig. 15), only our measure considers this cue and provides the best performance compared to the other features. Tab. 3 indicates that the pair-wise co-occurrence texture is suitable for FSS evaluation.

## 6 CONCLUSION

In this paper, we summarized the widely-used image quality measures for face sketch synthesis evaluation based on "pixel-level" degradation, and enumerated their limitations. Then we focused on the basic principles of sketching and proposed the use of structure co-occurrence texture as an alternative motivating principle for the design of face sketch style similarity measures. To amend the



**Figure 15: A toy example illustrating the difference between our pair-wise measure (right) and a pixel-wise measure (left).**

flaws of existing measures, we proposed the novel **Scoot-measure** which simultaneously captures both "block-level" spatial structure and co-occurrence texture. Also, we concluded that the pair-wise co-occurrence properties are more important than pixel-level properties for style similarity evaluation of face sketch synthesis images.

In addition, extensive experiments with 4 new meta-measures show that the proposed measure provides a reliable evaluation and achieves the best performance. Finally, we created two new human ranked face sketch datasets (containing 776 images & 1888 images) to examine the correlation between evaluation measures and human judgment. The study indicated a higher degree (78.8%) of correlation between our measure and human judgment. The proposed Scoot-measure is motivated from substantially different design principles, and we see it as complementary to the existing approaches.

All in all, our measure provides new insights into face sketch synthesis evaluation where current measures fail to truly examine the pros and cons of synthesis algorithms. We encourage the community to consider this measure in future model evaluations and comparisons.

**Future works.** Though our measure is only tested on face sketch datasets, it is not limited to evaluating face sketch synthesis. We will also investigate the potential to use the measure to describe styles in sketches and investigate new face sketch synthesis models. To promote the development of this field, our code and dataset will be made publicly available.

## REFERENCES

[1] Andrea Baraldi and Flavio Parmiggiani. 1995. An investigation of the textural characteristics associated with gray level cooccurrence matrix statistical parameters. *IEEE T Geosci. Remote.* 33, 2 (1995), 293–304.

[2] Itamar Berger, Ariel Shamir, Moshe Mahler, Elizabeth Carter, and Jessica Hodgins. 2013. Style and abstraction in portrait sketching. *ACM TOG* 32, 4 (2013), 55.

[3] DJ Best and DE Roberts. 1975. Algorithm AS 89: the upper tail probabilities of Spearman's rho. *J R STAT SOC C-APPL* 24, 3 (1975), 377–379.

[4] Lacey Best-Rowden, Hu Han, Charles Otto, Brendan F Klare, and Anil K Jain. 2014. Unconstrained face recognition: Identifying a person of interest from a media collection. *IEEE Transactions on Information Forensics and Security* 9, 12 (2014), 2144–2157.

[5] John Canny. 1987. A computational approach to edge detection. In *Readings in Computer Vision*. Elsevier, 184–203.

[6] David A Clausi. 2002. An analysis of co-occurrence texture statistics as a function of grey level quantization. *CAN J Remote. Sens.* 28, 1 (2002), 45–62.

[7] Deng-Ping Fan, Ming-Ming Cheng, Yun Liu, Tao Li, and Ali Borji. 2017. Structure-measure: A New Way to Evaluate Foreground Maps. *ICCV* (2017).

[8] William T Freeman, Joshua B Tenenbaum, and Egon C Pasztor. 2003. Learning style translation for the lines of a drawing. *ACM TOG* 22, 1 (2003), 33–46.

[9] Dennis Gabor. 1946. Theory of communication. Part 1: The analysis of information. *Journal of the Institution of Electrical Engineers-Part III: Radio and Communication Engineering* 93, 26 (1946), 429–441.

[10] Mary M Galloway. 1974. Texture analysis using grey level run lengths. *NASA STI/Recon Technical Report N* 75 (1974).

[11] Xinbo Gao, Juanjuan Zhong, Dacheng Tao, and Xuelong Li. 2008. Local face sketch synthesis learning. *Neurocomputing* 71, 10-12 (2008), 1921–1930.

[12] Stéphane Grabli, Emmanuel Turquin, Frédo Durand, and François X Sillion. 2004. Programmable style for NPR line drawing. *Rendering Techniques (Eurographics Symposium on Rendering)* (2004).

[13] Robert M Haralick. 1979. Statistical and structural approaches to texture. *Proc. IEEE* 67, 5 (1979), 786–804.

[14] Robert M Haralick, Karthikeyan Shanmugam, et al. 1973. Textural features for image classification. *IEEE Transactions on systems, man, and cybernetics* 6 (1973), 610–621.

[15] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. 2017. Image-To-Image Translation With Conditional Adversarial Networks. In *CVPR*. 1125–1134.

[16] Qingshan Liu, Xiaoou Tang, Hongliang Jin, Hanqing Lu, and Songde Ma. 2005. A nonlinear approach for face sketch synthesis and recognition. In *CVPR*, Vol. 1. IEEE, 1005–1010.

[17] Aude Oliva and Antonio Torralba. 2001. Modeling the shape of the scene: A holistic representation of the spatial envelope. *International journal of computer vision* 42, 3 (2001), 145–175.

[18] Chunlei Peng, Xinbo Gao, Nannan Wang, Dacheng Tao, Xuelong Li, and Jie Li. 2016. Multiple representations-based face sketch–photo synthesis. *IEEE TNNLS* 27, 11 (2016), 2201–2215.

[19] P Jonathon Phillips, Hyeonjoon Moon, Syed A Rizvi, and Patrick J Rauss. 2000. The FERET evaluation methodology for face-recognition algorithms. *IEEE TPAMI* 22, 10 (2000), 1090–1104.

[20] Jordi Pont-Tuset and Ferran Marques. 2013. Measures and meta-measures for the supervised evaluation of image segmentation. In *CVPR*. 2131–2138.

[21] Hamid R Sheikh and Alan C Bovik. 2006. Image information and visual quality. *IEEE TIP* 15, 2 (2006), 430–444.

[22] Irvin Sobel. 1990. An isotropic 3× 3 image gradient operator. *Machine vision for three-dimensional scenes* (1990), 376–379.

[23] Yibing Song, Linchao Bao, Qingxiong Yang, and Ming-Hsuan Yang. 2014. Real-time exemplar-based face sketch synthesis. In *ECCV*. Springer, 800–813.

[24] Bomba Sylwia, Cai Rovina, Croes Brun, Gerard Justin, and Lewis Marisa. 2015. *Beginner's Guide to Sketching*. 3dtotal Publishing.

[25] Xiaoou Tang and Xiaogang Wang. 2004. Face sketch recognition. *IEEE Transactions on Circuits and Systems for Video Technology* 14, 1 (2004), 50–57.

[26] Nannan Wang, Xinbo Gao, and Jie Li. 2018. Random sampling for fast face sketch synthesis. *Pattern Recognition* 76 (2018), 215–227.

[27] Nannan Wang, Xinbo Gao, Jie Li, Bin Song, and Zan Li. 2016. Evaluation on synthesized face sketches. *Neurocomputing* 214 (2016), 991–1000.

[28] Nannan Wang, Xinbo Gao, Leiyu Sun, and Jie Li. 2017. Bayesian face sketch synthesis. *IEEE Transactions on Image Processing* 26, 3 (2017), 1264–1274.

[29] Nannan Wang, Dacheng Tao, Xinbo Gao, Xuelong Li, and Jie Li. 2014. A comprehensive survey to face hallucination. *IJCV* 106, 1 (2014), 9–30.

[30] Nannan Wang, Shengchuan Zhang, Xinbo Gao, Jie Li, Bin Song, and Zan Li. 2017. Unified framework for face sketch synthesis. *Signal Processing* 130 (2017), 1–11.

[31] Nannan Wang, Mingrui Zhu, Jie Li, Bin Song, and Zan Li. 2017. Data-driven vs. model-driven: Fast face sketch synthesis. *Neurocomputing* (2017).

[32] Xiaogang Wang and Xiaoou Tang. 2009. Face photo-sketch synthesis and recognition. *IEEE TPAMI* 31, 11 (2009), 1955–1967.

[33] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE TIP* 13, 4 (2004), 600–612.

[34] Georges Winkenbach and David H Salesin. 1994. Computer-generated pen-and-ink illustration. In *ACM SIGGRAPH*. 91–100.

[35] Wufeng Xue, Lei Zhang, Xuanqin Mou, and Alan C Bovik. 2014. Gradient magnitude similarity deviation: A highly efficient perceptual image quality index.

[36] *IEEE Transactions on Image Processing* 23, 2 (2014), 684–695.

[36] Jun Yu, Meng Wang, and Dacheng Tao. 2012. Semisupervised multiview distance metric learning for cartoon synthesis. *IEEE Transactions on Image Processing* 21, 11 (2012), 4636–4648.

[37] Liliang Zhang, Liang Lin, Xian Wu, Shengyong Ding, and Lei Zhang. 2015. End-to-end photo-sketch generation via fully convolutional representation learning. In *ICMR*. ACM, 627–634.

[38] Lin Zhang, Lei Zhang, Xuanqin Mou, and David Zhang. 2011. FSIM: A feature similarity index for image quality assessment. *IEEE TIP* 20, 8 (2011), 2378–2386.

[39] Mingjin Zhang, Jie Li, Nannan Wang, and Xinbo Gao. 2017. Compositional Model-Based Sketch Generator in Facial Entertainment. *IEEE Transactions on Cybernetics* (2017).

[40] Shengchuan Zhang, Xinbo Gao, Nannan Wang, and Jie Li. 2016. Robust face sketch style synthesis. *IEEE Transactions on Image Processing* 25, 1 (2016), 220–232.

[41] Shengchuan Zhang, Xinbo Gao, Nannan Wang, and Jie Li. 2017. Face sketch synthesis from a single photo–sketch pair. *IEEE Transactions on Circuits and Systems for Video Technology* 27, 2 (2017), 275–287.

[42] Shengchuan Zhang, Xinbo Gao, Nannan Wang, Jie Li, and Mingjin Zhang. 2015. Face sketch synthesis via sparse representation-based greedy search. *IEEE transactions on image processing* 24, 8 (2015), 2466–2477.

[43] Wei Zhang, Xiaogang Wang, and Xiaoou Tang. 2011. Coupled information-theoretic encoding for face photo-sketch recognition. In *CVPR*. IEEE, 513–520.

[44] Yuqian Zhang, Nannan Wang, Shengchuan Zhang, Jie Li, and Xinbo Gao. 2016. Fast face sketch synthesis via kd-tree search. In *European Conference on Computer Vision*. Springer, 64–77.

[45] Hao Zhou, Zhanghui Kuang, and Kwan-Yee K Wong. 2012. Markov weight fields for face sketch synthesis. In *CVPR*. IEEE, 1091–1097.

[46] Mingrui Zhu, Nannan Wang, Xinbo Gao, and Jie Li. 2017. Deep graphical feature learning for face sketch synthesis. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*. AAAI Press, 3574–3580.

[47] Steven W Zucker. 2005. Which computation runs in visual cortical columns? *23 Problems in Systems Neuroscience* (2005), 215.

[48] Steven W Zucker, Allan Dobbins, and Lee Iverson. 1989. Two stages of curve detection suggest two styles of visual computation. *Neural computation* 1, 1 (1989), 68–81.