

Scoot: A Perceptual Metric for Facial Sketches

Supplemental Material

Deng-Ping Fan^{1,2} ShengChuan Zhang³ Yu-Huan Wu¹ Yun Liu¹

Ming-Ming Cheng^{1,*} Bo Ren¹ Paul L. Rosin⁴ Rongrong Ji³

¹ TKLNDST, CS, Nankai University ² IIAI ³ Xiamen University ⁴ Cardiff University

Abstract

In this document, we include additional materials related to the benchmark, proposed dataset, and results. All of the benchmark results will be released in our website.

- **Benchmark.** We systematically assess 9 representative measures [1–6, 8–10] over two (CUFS [7] and CUFSF [11]) widely-used datasets, making it the first one largest-scale measure benchmark in this field. Extensive experiments on this study verify that our Scoot measure exceeds performance of prior works.
- **Dataset.** We collect two new human-ranked datasets, i.e., RCUFS and RCUFSF (3.8k face sketches in total), making it the largest-scale subjectively verified datasets.
- **Results.** We provide visualization results on the proposed 3 meta-measures and the proposed human-ranked datasets.

1. Benchmark

As shown in Fig. 1, it has witnessed the dramatic development of face sketch modeling, while the community long-term suffered from the lack of a standard representative perceptual metric based benchmark. To the best of our knowledge, this work is the first and the largest-scale benchmark for this issue.

2. Dataset

In Fig. 2, we present the example of the proposed human-ranked dataset. We use these datasets to examine the ranking consistency between current measures and human judgments. The detail can be found in Section 5.3 (manuscript). We refer the reader to the accompanying attachment (“Proposed Datasets”) for the final human-ranked datasets.

*M.M. Cheng (cmm@nankai.edu.cn) is the corresponding author.

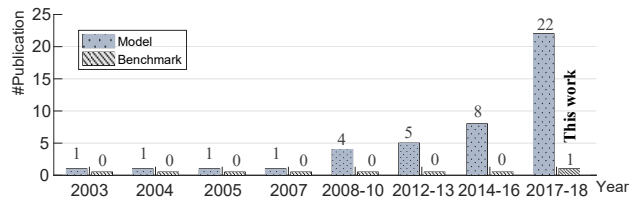


Figure 1: Number of FSS papers on top conferences (i.e., CVPR, ICCV, ECCV), ACM TOG, and IEEE Trans. journals over the past 16 years.

3. Results

3.1. Meta-measure 1: Stability to Slight Resizing.

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

Fig. 7 and Fig. 8 show the results: the narrower the band is, the more stable a measure is to slightly downsizing. We can see our Scoot measure achieves a significant improvement over the existing measures, such as GMS-D [9], SSIM [8], FSIM [10], and VIF [4] measures in both the CUFS and CUFSF databases. Among the alternative texture measures, our measure is the best. These improvements are mainly because the proposed measure considers “block-level” statistics rather than “pixel-level”.

One example shows in Fig. 3 which illustrates our Scoot measure is more robust than existing widely-used measures. For SSIM measure, it switches the ranking of LR and SSD when using GT or Resized-GT as reference. For other measures (VIF, FSIM), they also change the ranking in different ways. Only our measure keeps the original ranking.

3.2. Meta-measure 2: Rotation Sensitivity.

In real-world situations, sketches drawn by artists may also have slight rotations compared to the original photographs. Thus, the proposed second meta-measure verifies the sensitivity of GT rotation for the evaluation measure. We did a slight counter-clockwise rotation (5°) for each GT sketch.



Figure 2: 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. We refer the reader to the accompanying attachment (“*Proposed Datasets*”) for more details.

The sensitivity results are shown in Fig. 7 and Fig. 8. The thinner the band is, the better the measure performs. Our measure also significantly outperforms the current measures over the CUFS and CUFSF databases. We attribute the robust performance to the exploit of local region-level statistics rather than pixel-level statistics.

We show an example in Fig. 4. For SIM, VIF and FSIM, they change the ranking result. However, our measure correctly ranks these results and does not change their ranking when using different GT as reference.

3.3. Meta-measure 3: Content Capture Capability.

The third meta-measure describes that a good measure should assign a complete sketch generated by SOTA algorithm higher score than the sketches of only preserving incomplete strokes.

From Fig. 9 and Fig. 10, we observe a great improvement over the other measures in CUFS database. A slight improvement is also achieved for the CUFSF database.

Fig. 5 shows that our measure assigns the complete sketch generated by state-of-the-art algorithm (e.g., MRF [7]) higher score than the sketch of only preserving light (*incomplete*) strokes.

3.4. Meta-measure 4: Human Judgment.

The fourth meta-measure (Jug) specifies that the ranking result according to an evaluation measure should agree with the human judgment.

For human judgment in Fig. 9 and Fig. 10, the proposed Scoot measure shows a great 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.

Various of examples in Fig. 6 demonstrate that our measure provide an reliable evaluation.

References

- [1] J. Canny. A computational approach to edge detection. In *Readings in Computer Vision*, pages 184–203. Elsevier, 1987. 1, 3, 7, 8
- [2] D. Gabor. Theory of communication. part 1: The analysis of information. *Journal of the Institution of Electrical Engineers-Part III: Radio and Communication Engineering*, 93(26):429–441, 1946. 1, 3, 7, 8
- [3] M. M. Galloway. Texture analysis using grey level run lengths. *NASA STI/Recon Technical Report N, 75*, 1974. 1, 3, 7, 8
- [4] H. R. Sheikh and A. C. Bovik. Image information and visual quality. *IEEE TIP*, 15(2):430–444, 2006. 1, 3, 4, 5, 6, 7
- [5] H. R. Sheikh, A. C. Bovik, and G. De Veciana. An information fidelity criterion for image quality assessment using natural scene statistics. *IEEE TIP*, 14(12):2117–2128, 2005. 1, 3
- [6] I. Sobel. An isotropic 3×3 image gradient operator. *Machine vision for three-dimensional scenes*, pages 376–379, 1990. 1, 3, 7, 8

Dataset	Meta-measure	Classical Measure					Scoot	Improvement
		IFC [5]	VIF [4]	GMSD [9]	FSIM [10]	SSIM [8]		
CUFS	MM1↓ ranking	0.256 3	0.322 5	0.417 6	0.268 4	0.162 2	0.037 1	+77.2% = ((0.162 - 0.037)/0.162)
	MM2↓ ranking	0.189 4	0.236 6	0.210 5	0.123 3	0.086 2	0.025 1	+70.9% = ((0.086 - 0.025)/0.086)
	MM3↑ ranking	1.20% 6	43.5% 3	21.9% 4	14.2% 5	81.4% 2	95.9% 1	+14.5% = (95.9% - 81.4%)
	Jud↑ ranking	26.9% 6	44.1% 3	42.6% 4	50.0% 2	37.3% 5	76.3% 1	+26.3% = (76.3% - 50.0%)
	overall ranking	6	5	4	3	2	1	
	CUFSF	MM1↓ ranking	0.089 3	0.111 4	0.259 6	0.151 5	0.073 2	0.012 1
MM2↓ ranking		0.112 4	0.150 6	0.132 5	0.058 2	0.074 3	0.008 1	+86.2% = ((0.058 - 0.008)/0.058)
MM3↑ ranking		3.07% 6	22.2% 5	63.6% 3	32.4% 4	97.4% 2	97.5% 1	+0.10% = (97.5% - 97.4%)
Jud↑ ranking		25.4% 6	52.8% 3	58.6% 2	37.5% 4	36.8% 5	78.8% 1	+20.2% = (78.8% - 80.9%)
overall ranking		6	5	4	3	2	1	
		Alternative Texture Measure						
Dataset	Meta-measure	Canny [1]	GLRLM [3]	Sobel [6]	Gabor [2]	Scoot	Improvement	
CUFS	MM1↓ ranking	0.086 4	0.111 5	0.040 2	0.062 3	0.037 1	+7.50% = ((0.04 - 0.037)/0.04)	
	MM2↓ ranking	0.078 4	0.111 5	0.037 2	0.055 3	0.025 1	+32.4% = ((0.037 - 0.025)/0.037)	
	MM3↑ ranking	33.7% 2	18.6% 3	0.00% 4	0.00% 4	95.9% 1	+62.2% = (95.9% - 33.7%)	
	Jud↑ ranking	27.8% 5	73.7% 2	32.8% 4	72.2% 3	76.3% 1	+2.60% = (76.3% - 73.7%)	
	overall ranking	5	4	3	2	1		
	CUFSF	MM1↓ ranking	0.138 5	0.125 4	0.048 2	0.089 3	0.012 1	+75.0% = ((0.048 - 0.012)/0.048)
MM2↓ ranking		0.146 5	0.079 4	0.044 3	0.043 2	0.008 1	+81.4% = ((0.043 - 0.008)/0.043)	
MM3↑ ranking		0.00% 4	64.6% 2	0.00% 4	19.3% 3	97.5% 1	+32.9% = (97.5% - 64.6%)	
Jud↑ ranking		0.10% 5	68.0% 3	52.6% 4	80.9% 1	78.8% 2	-2.10% = (78.8% - 80.9%)	
overall ranking		5	4	3	2	1		

Table 1: Benchmarking results of 5 classical measures and 4 alternative textures on the CUFS and CUFSF datasets. Darker color indicates better performance. These differences are all statistically significant at the $\alpha < 0.05$ level.

[7] X. Wang and X. Tang. Face photo-sketch synthesis and recognition. *IEEE TPAMI*, 31(11):1955–1967, 2009. 1, 2, 6

[8] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE TIP*, 13(4):600–612, 2004. 1, 3, 4, 5, 6, 7

[9] W. Xue, L. Zhang, X. Mou, and A. C. Bovik. Gradient magnitude similarity deviation: A highly efficient perceptual image quality index. *IEEE TIP*, 23(2):684–695, 2014. 1, 3, 7

[10] L. Zhang, L. Zhang, X. Mou, and D. Zhang. FSIM: A feature similarity index for image quality assessment. *IEEE TIP*, 20(8):2378–2386, 2011. 1, 3, 4, 5, 6, 7

[11] W. Zhang, X. Wang, and X. Tang. Coupled information-theoretic encoding for face photo-sketch recognition. In *IEEE CVPR*, pages 513–520, 2011. 1

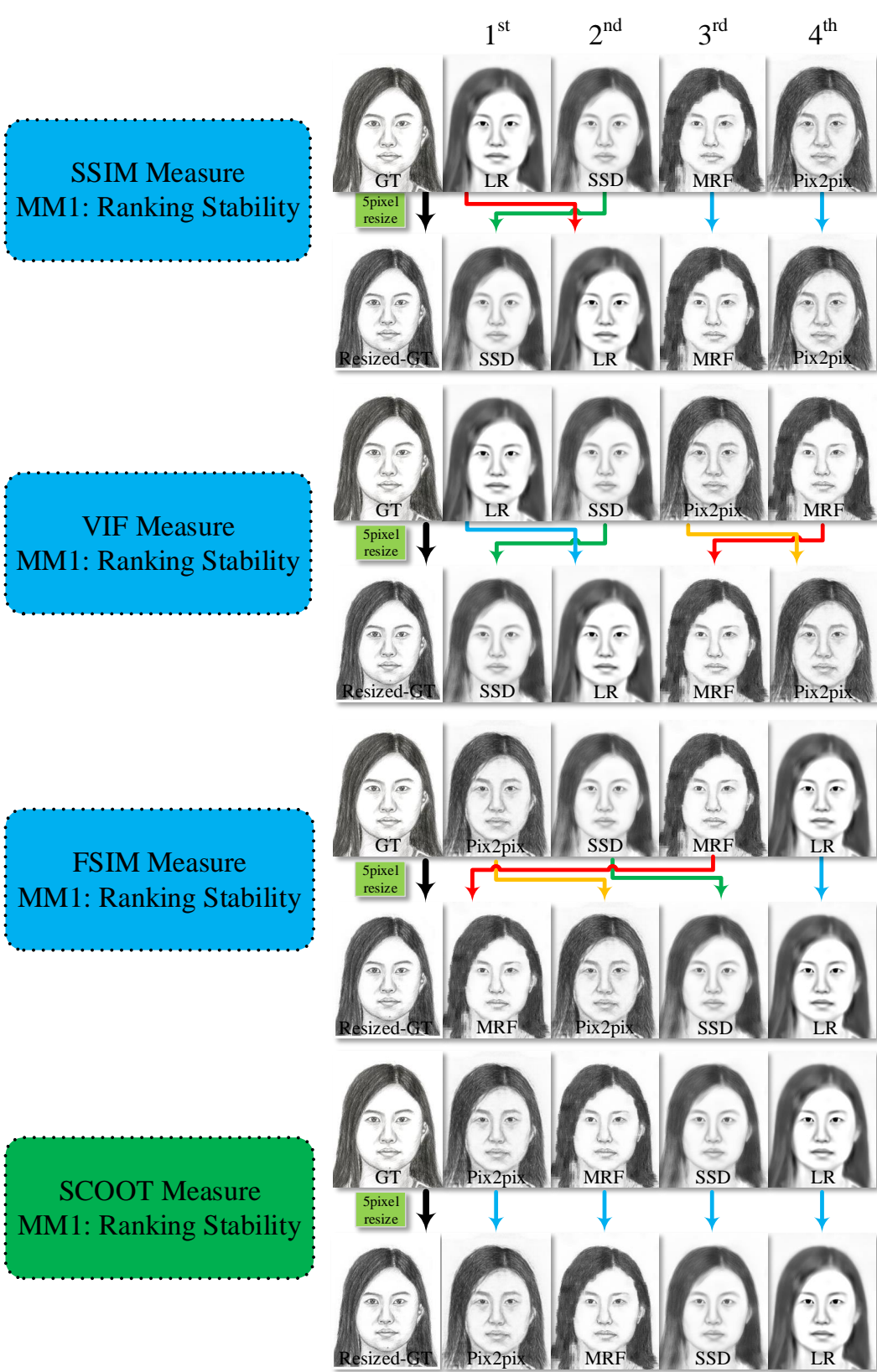


Figure 3: Visual comparison of existing widely-used FSS measures (SSIM [8], FSIM [10], and VIF [4]) on meta-measure 1. The experiment clearly shows that the proposed SCOOT measure is more stable to slightly resizing.

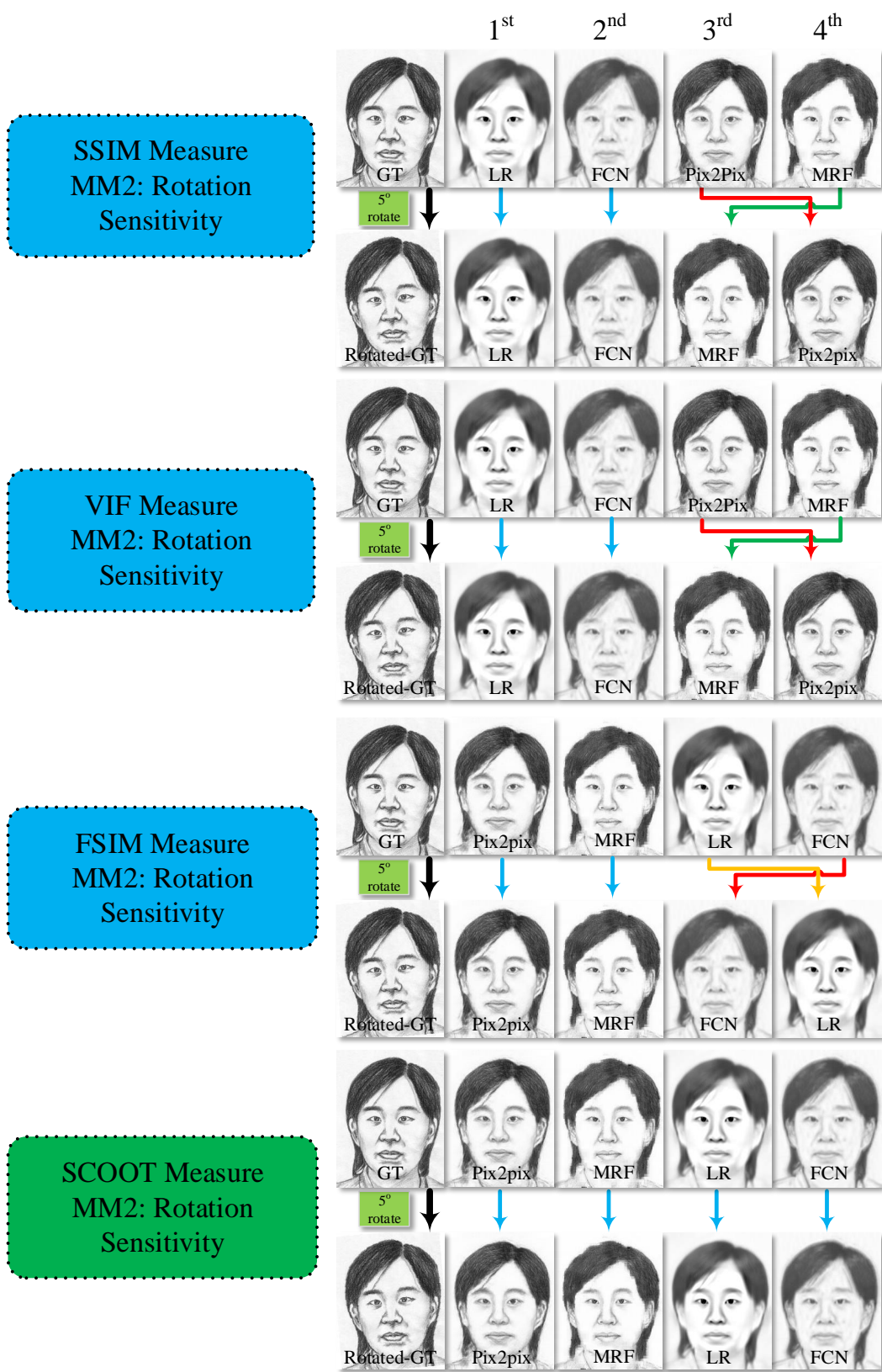


Figure 4: Visual comparison of existing widely-used FSS measures (SSIM [8], FSIM [10], and VIF [4]) on meta-measure 2. The experiment clearly demonstrates that the proposed SCOOT measure is less sensitive to minor rotation.

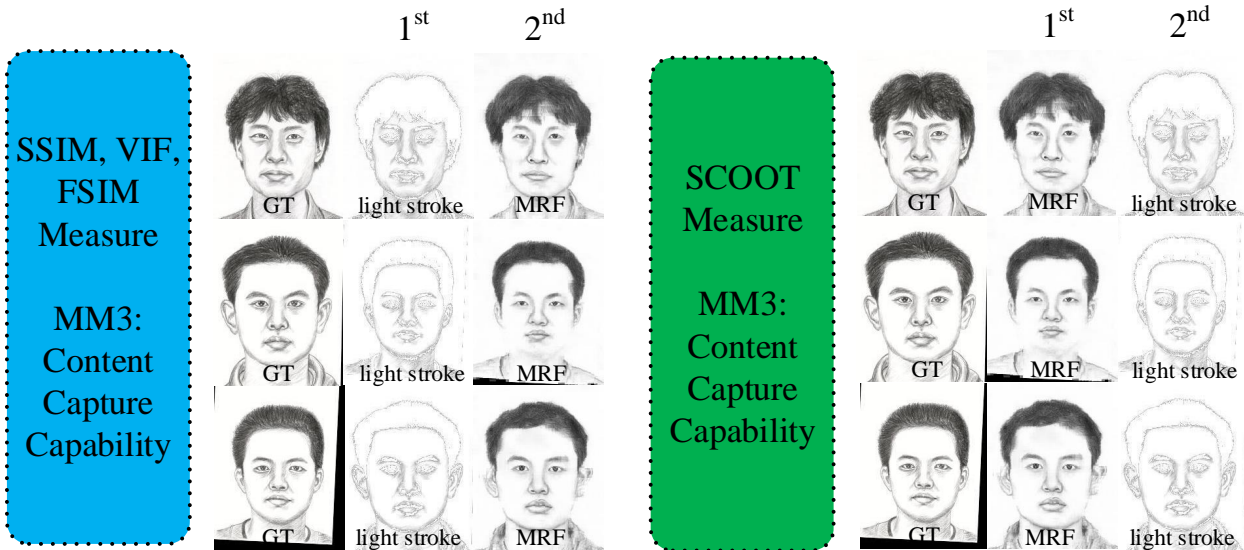


Figure 5: Visual comparison of existing widely-used FSS measures (SSIM [8], FSIM [10], and VIF [4]) on meta-measure 3. We conduct this experiment to verify the content capture capability of the measures. Our measure assigns the complete sketch generated by state-of-the-art algorithm (e.g., MRF [7]) higher score than the sketch of only preserving light (*incomplete*) strokes.

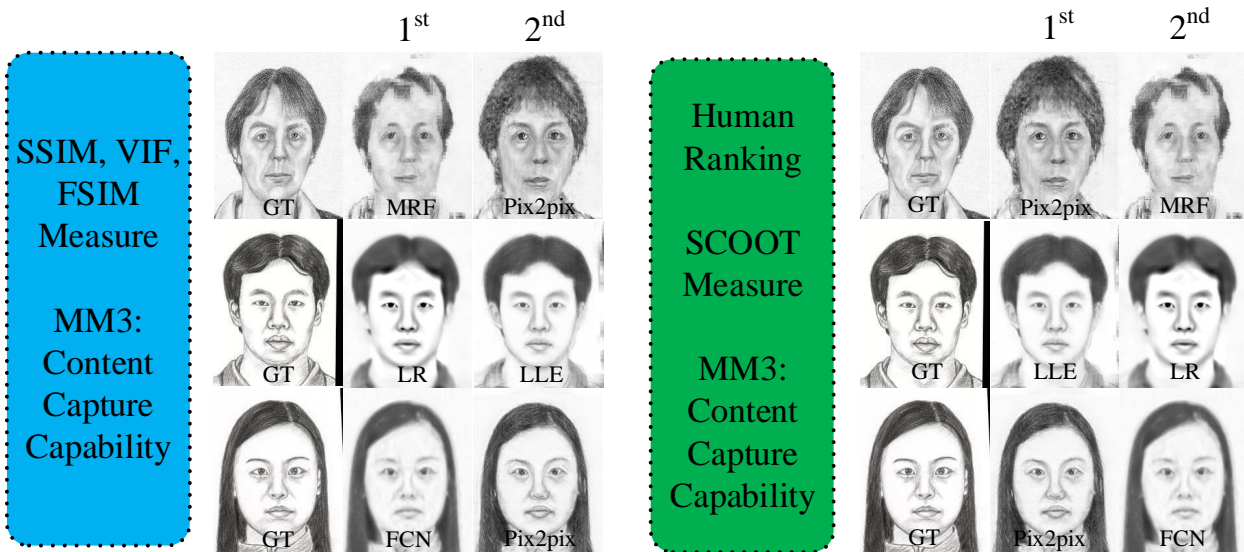


Figure 6: Visual comparison of existing popular FSS measures (SSIM [8], FSIM [10], and VIF [4]) on meta-measure 4. The proposed Scoot measure shows higher consistency with human judgment than previous methods on our collected datasets. Zoom-in for details.

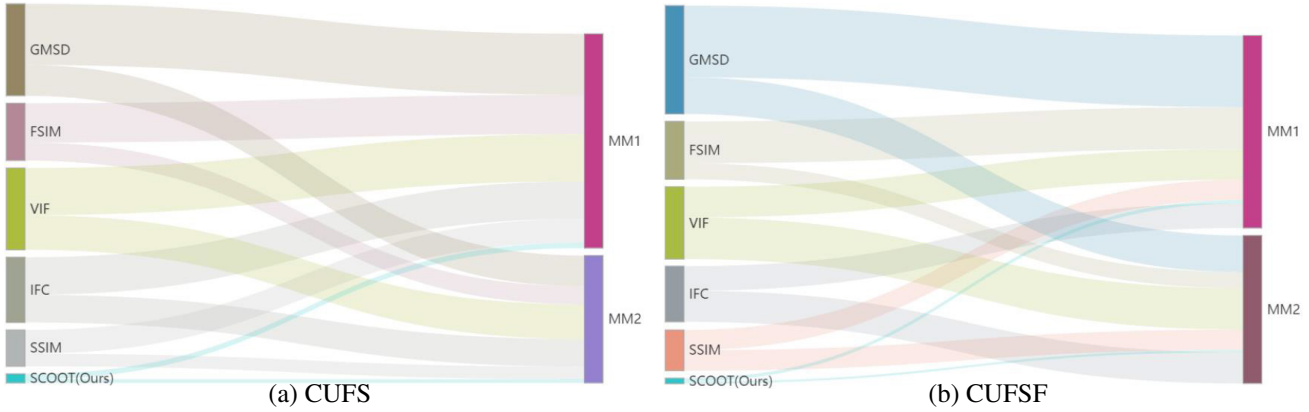


Figure 7: **Quantitative comparison of different measure (GMSD [9], SSIM [8], FSIM [10], VIF [4], and the proposed Scoot) on our meta-measure 1 & 2.** The narrower the color band is, the better the performance is. The proposed Scoot measure (in the last row) achieves the best performance. Zoom-in for details.

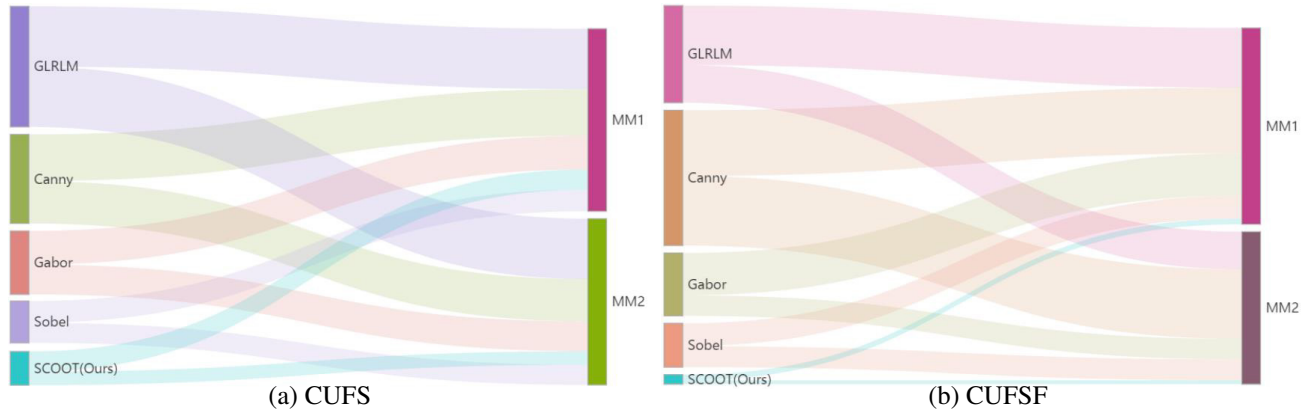


Figure 8: **Quantitative comparison of alternative texture (Canny [1], Sobel [6], GLRLM [3], Gabor [2] and the proposed Scoot) on our meta-measure 1 & 2.** The narrower the color band is, the better the performance is. Our Scoot measure (in the last row) achieves the preferable performance.

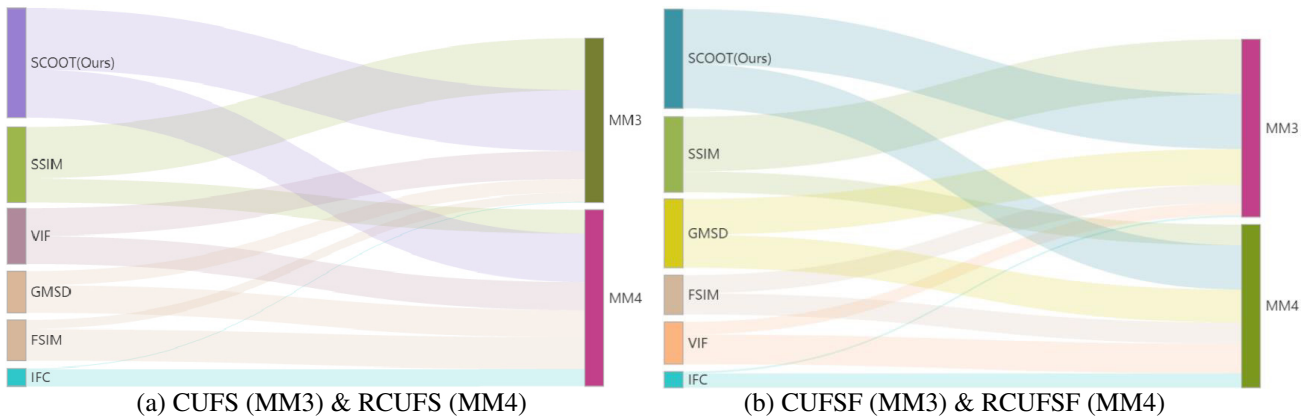


Figure 9: **Quantitative comparison of different measure (GMSD [9], SSIM [8], FSIM [10], VIF [4], and the proposed Scoot) on our meta-measure 3 & 4.** The wider the color band is, the better the performance is. Our Scoot measure (1st row) outperforms prior works. Zoom-in for details.

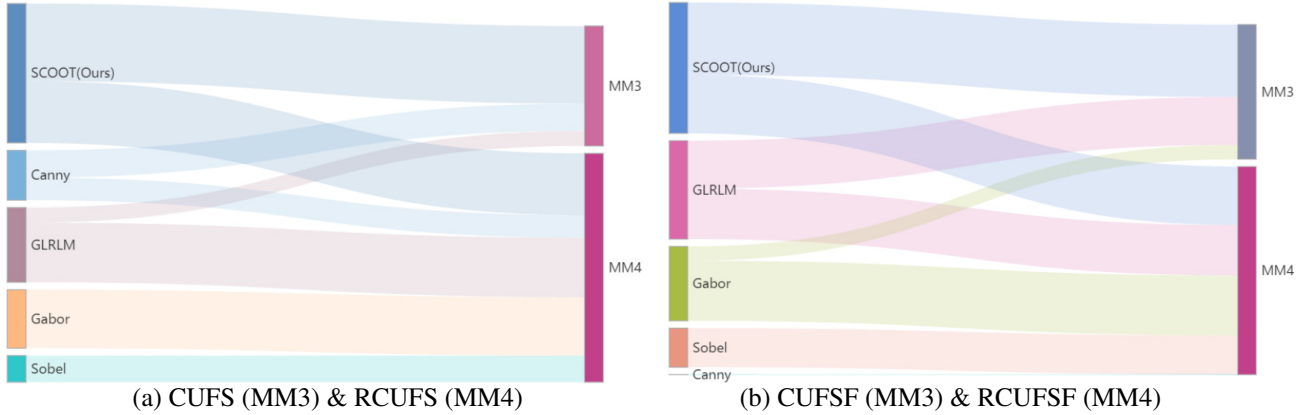


Figure 10: **Quantitative comparison of alternative texture (Canny [1], Sobel [6], GLRLM [3], Gabor [2] and the proposed Scoot) on our meta-measure 3 & 4.** The wider the color band is, the better the performance is. Our measure present in first row provides the best performance in both MM3 and MM4.