

# Recognition of Action Units in the Wild with Deep Nets and a New Global-Local Loss

C. Fabian Benitez-Quiroz

Yan Wang

Alex M. Martinez

Dept. Electrical and Computer Engineering  
The Ohio State University

{benitez-quiroy.1,wang.9021,martinez.158}@osu.edu

## 1. Derivations of Eq. (9)

The main results given in (9) in the main paper was derived using Green's Theorem. Recall, that Green's Theorem is a mechanism to convert a line integral about a Jordan curve  $\mathcal{C}$  to a double integral over the plane  $\mathcal{D}$  bounded by  $\mathcal{C}$ . The Jordan curve  $\mathcal{C}$  must be positive oriented and piecewise smooth.

We begin by defining the area of  $\mathcal{D}$  we wish to compute. We assume we are in  $\mathbb{R}^2$ , defined by the  $(x_1, x_2)$  coordinate system. Using Green's Theorem, we have

$$\text{Area}(\mathcal{D}) = \frac{1}{2} \oint_{\mathcal{C}} -x_2 dx_1 + x_1 dx_2, \quad (\text{S1})$$

where  $\mathcal{D}$  is a non-self-intersecting polygon bounded by  $\mathcal{C}$ .

In our case, this polygon and bounding Jordan curve  $\mathcal{C}$  are defined by a set of  $t$  points on  $\mathcal{C}$ . An example of this curve was shown in Figure 3 (right-most image) in the main paper. These points were defined as  $(\tilde{x}_{i1}, \dots, \tilde{x}_{it})$ , where  $\tilde{x}_{ik} = (\tilde{x}_{ik1}, \tilde{x}_{ik2})^T$ .

We call the line segment from each  $\tilde{x}_{ik}$  to  $\tilde{x}_{i(k+1)}$ ,  $\Gamma_i$ . A parametrization of  $\Gamma_i$  is given by the function  $\gamma_i : [0, 1] \rightarrow \mathbb{R}^2$  and, hence, we have  $\gamma_i(a) = (\tilde{x}_{ik1}, \tilde{x}_{ik2}) + a(\tilde{x}_{i(k+1)1} - \tilde{x}_{ik1}, \tilde{x}_{i(k+1)2} - \tilde{x}_{ik2})$ .

We can now compute the line integral of this curve,

$$\begin{aligned} \oint_{\Gamma_i} -x_2 dx_1 + x_1 dx_2 &= \\ & \int_0^1 (-\tilde{x}_{ik2} - a(\tilde{x}_{i(k+1)2} - \tilde{x}_{ik2}))(\tilde{x}_{i(k+1)1} - \tilde{x}_{ik1}) da \\ & + \int_0^1 (-\tilde{x}_{ik1} - a(\tilde{x}_{i(k+1)1} - \tilde{x}_{ik1}))(\tilde{x}_{i(k+1)2} - \tilde{x}_{ik2}) da \\ & = \int_0^1 (-\tilde{x}_{ik2}\tilde{x}_{i(k+1)1} + \tilde{x}_{ik1}\tilde{x}_{i(k+1)2}) da \\ & = \tilde{x}_{ik1}\tilde{x}_{i(k+1)2} - \tilde{x}_{ik2}\tilde{x}_{i(k+1)1}. \end{aligned} \quad (\text{S2})$$

Hence, the area of  $\mathcal{D}$  is

$$\frac{1}{2} \left[ \left( \sum_{k=1}^{t-1} (\tilde{x}_{ik1}\tilde{x}_{i(k+1)2} - \tilde{x}_{ik2}\tilde{x}_{i(k+1)1}) \right) + (\tilde{x}_{it1}\tilde{x}_{i12} - \tilde{x}_{i12}\tilde{x}_{it1}) \right]. \quad (\text{S3})$$

## 2. Extended Experimental Results

The main paper presented comparative results with state-of-the-art algorithms using the *Final score* of the EmotionNet challenge, given in (18). Herein, we provide results with the two metric defining this final score:  $F_1$  score and *Accuracy*, given in (16) and (17), respectively. Figures S1, S2 and S3 plot the *Accuracy* values of the experiments described in Sections 4.1 and 4.2. Figures S4, S5 and S6 show the  $F_1$  scores of these same experiments. Tables S1, S2, S3 and S4 show the results from figures 7-9 in the main manuscript. A right tailed  $t$ -test shows that our method is statistically better than state-of-the-art algorithms ( $p < .0005$ ,  $p < .005$  and  $p < 10^{-8}$  with respect to JHU, I2R-CCNU-NTU-2 and AlexNet).

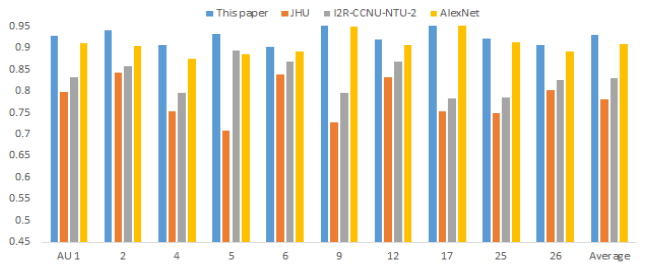


Figure S1. Results on the EmotionNet testing dataset. *Accuracy* calculated using (17).

## 3. Extended Experimental Results on Landmark detection

Comparative results are given against state-of-the-art algorithms and the top performers on the 300-W challenge

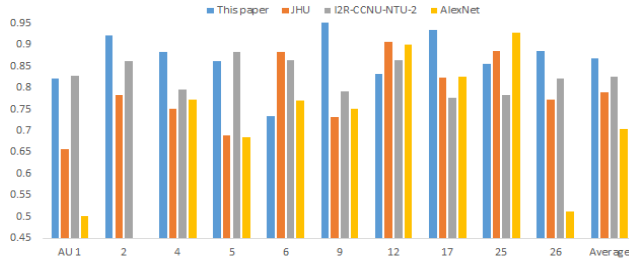


Figure S2. Average *Accuracy* for images at different scales.

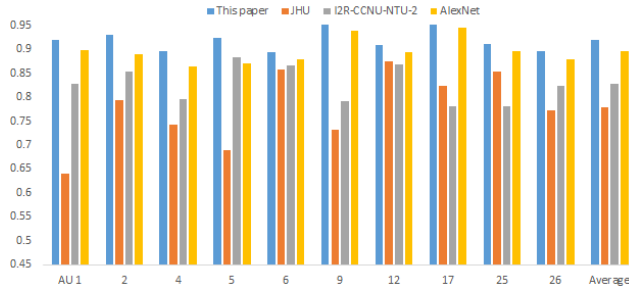


Figure S3. *Accuracy* for images with small occluders.

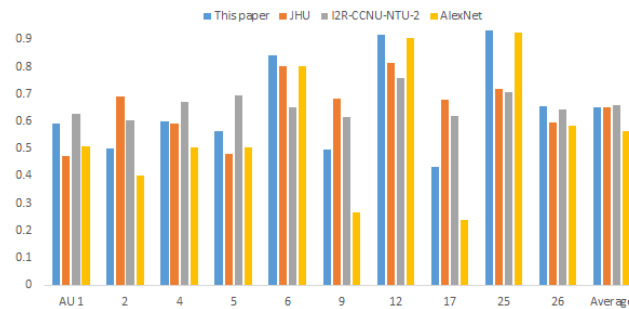


Figure S4. Results on the EmotionNet testing dataset.  $F_1$  score is calculated using (16).

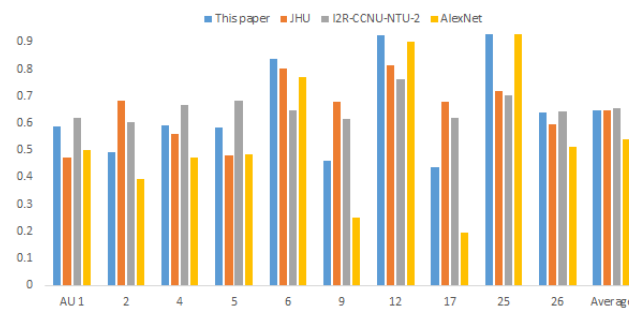


Figure S5. Average  $F_1$  score for images at different scales.

dataset [1]. This dataset includes a large number of image variations and a diverse group of people of distinct ethnic and cultural backgrounds. The images are divided following the protocol of [2]: 3,148 images are used for training (to which we apply our data augmentation approach defined above) and 689 faces serve as the testing set. We compare

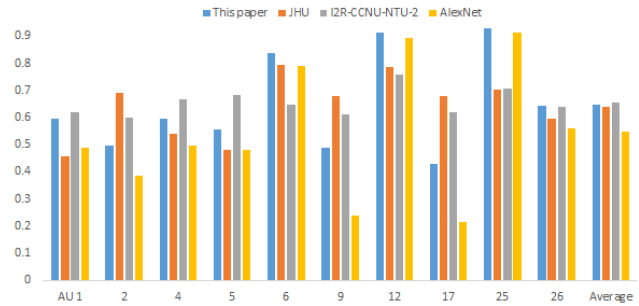


Figure S6.  $F_1$  scores for images with small occluders.

Table S1. Final scores on the EmotionNet Dataset (given by (18)).

AU	This paper	JHU	I2R-CCNU-NTU-2	AlexNet
1	<b>.76</b>	.57	.73	.71
2	.72	<b>.74</b>	<b>.73</b>	.65
4	<b>.75</b>	.67	.73	.69
5	.75	.58	<b>.79</b>	.7
6	<b>.87</b>	.84	.76	.85
9	<b>.73</b>	<b>.72</b>	.71	.61
12	<b>.92</b>	.86	.81	<b>.91</b>
17	.7	<b>.75</b>	.7	.6
25	<b>.93</b>	.8	.75	<b>.92</b>
26	<b>.78</b>	.69	.74	.74
Average	<b>.79</b>	.72	.74	.74

Table S2.  $F_1$  scores on the EmotionNet Dataset (given by (16)).

AU	This paper	JHU	I2R-CCNU-NTU-2	AlexNet
1	.59	.47	<b>.63</b>	.51
2	.5	<b>.69</b>	.6	.4
4	.6	.59	<b>.67</b>	.51
5	.56	.48	<b>.69</b>	.5
6	<b>.84</b>	.8	.65	.8
9	.5	<b>.68</b>	.62	.27
12	<b>.92</b>	.81	.76	.9
17	.43	<b>.68</b>	.62	.24
25	<b>.93</b>	.72	.71	<b>.93</b>
26	<b>.66</b>	.6	<b>.65</b>	.58
Average	<b>.65</b>	<b>.65</b>	<b>.66</b>	.56

our results with top performing algorithms: Explicit Shape Regression (ESR) [3], Supervised Descent Method (SDM) [4] and Local Binary Features (LBF) [2].

The detection error is the point-to-point squared Euclidean distance,  $\|\mathbf{f}_i - \mathbf{y}_i\|_2^2$ , normalized by the Euclidean distance between the outer corners of the eyes.

As shown in Table S5, our proposed global method achieves the smallest error. Additionally, Figure S7 compares the results of the proposed GL-CNN algorithm against the top performers in the 300-W challenge [5, 6]. As seen in

Table S3. Average scores (given by (16)) for images at different scales.

AU	This paper	JHU	I2R-CCNU-NTU-2	AlexNet
1	<b>.71</b>	.57	<b>.72</b>	.65
2	.71	<b>.73</b>	<b>.73</b>	.59
4	<b>.74</b>	.66	<b>.73</b>	.68
5	.72	.58	<b>.78</b>	.63
6	.79	<b>.84</b>	.76	.8
9	<b>.71</b>	.7	.7	.58
12	<b>.88</b>	.86	.82	<b>.88</b>
17	.69	<b>.75</b>	.7	.55
25	<b>.89</b>	.8	.74	.86
26	<b>.76</b>	.69	.73	.71
Average	<b>.76</b>	.72	.74	.69

Table S4. Final scores (equation (18)) for images with small occluders.

AU	This paper	JHU	I2R-CCNU-NTU-2	AlexNet
1	<b>.76</b>	.55	.72	.69
2	.71	<b>.74</b>	<b>.73</b>	.64
4	<b>.75</b>	.64	.73	.68
5	.74	.58	<b>.78</b>	.68
6	<b>.87</b>	.83	.76	.84
9	<b>.72</b>	.7	.7	.59
12	<b>.91</b>	.83	.81	.89
17	.69	<b>.75</b>	.7	.58
25	<b>.92</b>	.78	.75	<b>.91</b>
26	<b>.77</b>	.69	.73	.72
Average	<b>.78</b>	.71	.74	.72

Method	Mean normalized error
ESR	7.58
SDM	7.52
LBF	6.32
GL-CNN	<b>5.47</b>

Table S5. Average mean normalized error on the 300-W challenge dataset.

this plot, the proposed approach yields results comparable or slightly better than these algorithms but with the added advantage that our method can run at >60 frames/s in Matlab on a i7 desktop.

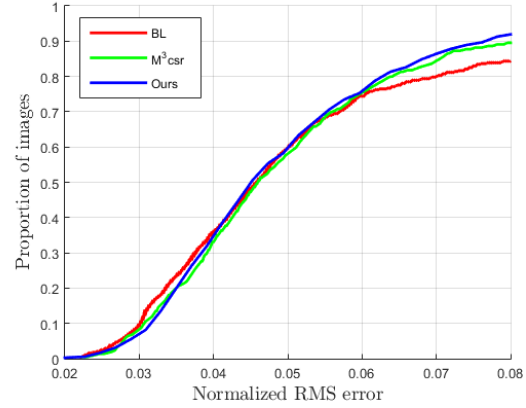


Figure S7. Cumulative normalized root mean square (RMS) error for the top performers on the 300-W challenge [5, 6] and our proposed GL-CNN algorithm. The  $y$ -axis specifies the proportion of images, with 1 indicating all images in the database are included. Note that the proposed algorithm outperforms the others when all images are included; demonstrating not only good local fits, but global (overall) fits as well.

## References

- [1] C. Sagonas, G. Tzimiropoulos, S. Zafeiriou, and M. Pantic, “A semi-automatic methodology for facial landmark annotation,” in *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR-W), 5th Workshop on Analysis and Modeling of Faces and Gestures*, 2013. 2
- [2] S. Ren, X. Cao, Y. Wei, and J. Sun, “Face alignment at 3000 fps via regressing local binary features,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1685–1692, 2014. 2
- [3] X. Cao, Y. Wei, F. Wen, and J. Sun, “Face alignment by explicit shape regression,” *International Journal of Computer Vision*, vol. 107, no. 2, pp. 177–190, 2014. 2
- [4] X. Xiong and F. De La Torre, “Supervised descent method and its applications to face alignment,” in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 532–539, 2013. 2
- [5] J. Deng, Q. Liu, J. Yang, and D. Tao, “M<sup>3</sup> csr: Multi-view, multi-scale and multi-component cascade shape regression,” *Image and Vision Computing*, 2015. 2, 3
- [6] H. Fan and E. Zhou, “Approaching human level facial landmark localization by deep learning,” *Image and Vision Computing*, 2015. 2, 3