

# A Deep Regression Architecture with Two-Stage Re-initialization for High Performance Facial Landmark Detection

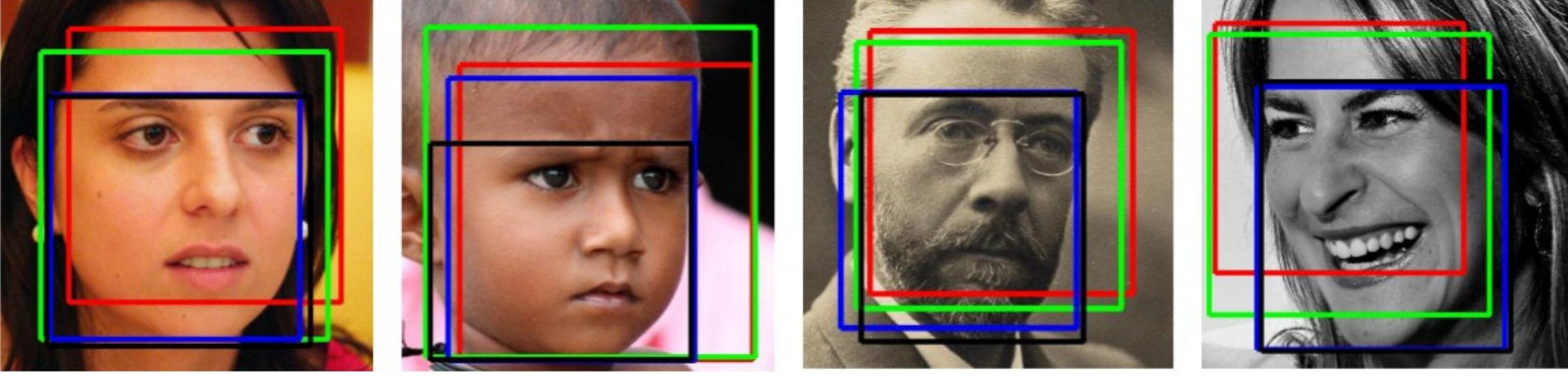
Jiangjing Lv<sup>1\*</sup>, Xiaohu Shao<sup>1\*</sup>, Junliang Xing<sup>2</sup>, Cheng Cheng<sup>1</sup>, Xi Zhou<sup>1</sup>

<sup>1</sup>Chongqing Institute of Green and Intelligent Technology, Chinese Academy of Sciences

<sup>2</sup>National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences

## Observations

- Most of facial landmark detection algorithms are sensitive to the initialization brought by face detection results.
- Different face detectors often return various face bounding boxes with different scales and center shifts.



## Motivations

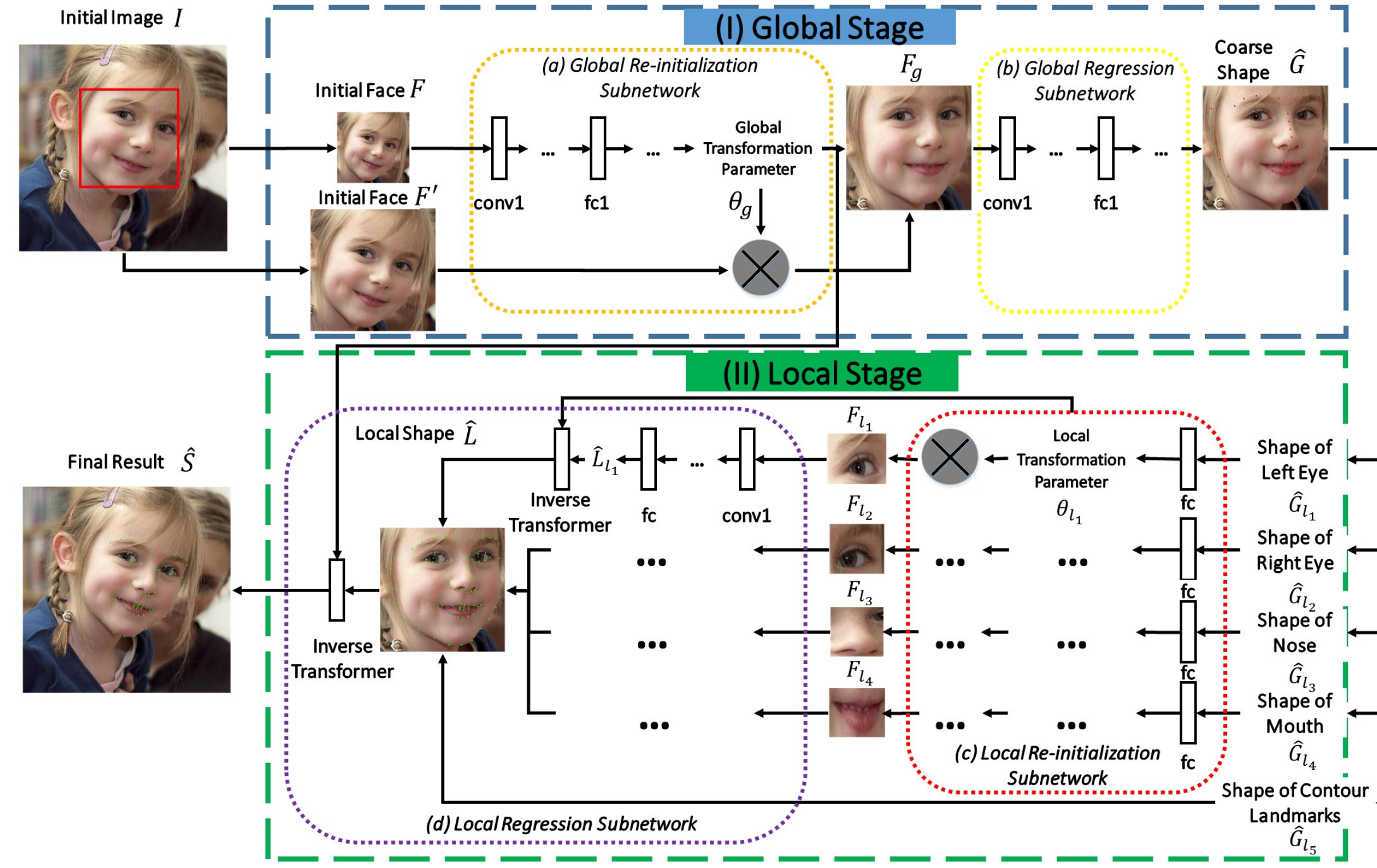
- Face boxes bounding by the ground truth shapes can boost the accuracy of facial landmark detection.
- Good performance of STN [14] on learning instance-specific transformations of training samples.

## Contributions

- An end-to-end deep regression with two-stage re-initialization architecture.
- Supervised spatial transformer learning.
- Good robustness to different initialization.
- State-of-the-art performances on 300-W and AFLW.

## Our Approach

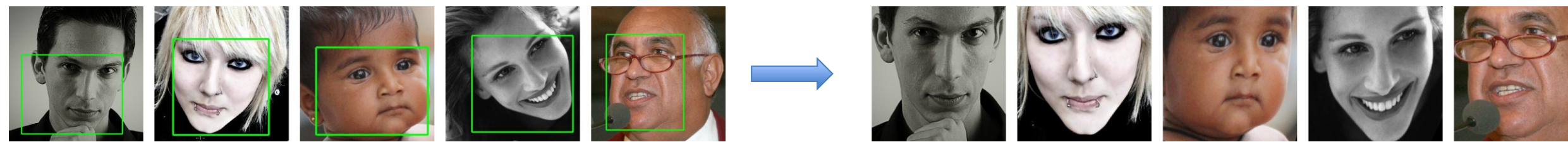
- Our two-stage re-initialization model successively re-initializes a deep regression model from coarse to fine, and global to local, to substantially boost the landmark detection performance.
- Global Stage:
  - Global Re-initialization Subnetwork
  - Global Regression Subnetwork
- Local Stage:
  - Local Re-initialization Subnetwork
  - Local Regression Subnetwork



## Global Stage

- The best canonical state for initialization of following regression is obtained by the supervised transformer subnetwork.

$$F_g = A(F', T_{\theta_g})$$



- The target of the global regression is the transformed ground truth shape:

$$G^* = T_{\theta_g}^{-1} \begin{pmatrix} S^* \\ 1 \end{pmatrix}$$

## Local Stage

- Different parts of the face shape are further separately re-initialized to their own canonical states.



- Local regression to minimize the loss function of shape increment.
- Inverse Transformer:  $\hat{S}_{l_n} = T_R T_{\theta_g} T_{\theta_{l_n}} \begin{pmatrix} \hat{I}_{l_n} \\ 1 \end{pmatrix}$

## Experiments

We first evaluate the robustness of our approach for various initialization, and then compare it with other state-of-the-art methods on the benchmark datasets.

Detectors	Common Subset	Challenging Subset	Full Set
VJ <sub>B<sub>1</sub></sub>	8.90	14.39	9.98
Dlib <sub>B<sub>1</sub></sub>	6.88	12.40	7.96
OD <sub>B<sub>1</sub></sub>	5.43	8.97	6.12
GT <sub>B<sub>1</sub></sub>	5.24	7.65	5.71
VJ <sub>B<sub>2</sub></sub>	6.19	10.15	6.96
Dlib <sub>B<sub>2</sub></sub>	5.30	9.13	6.05
OD <sub>B<sub>2</sub></sub>	5.03	8.43	5.69
GT <sub>B<sub>2</sub></sub>	5.04	7.64	5.55
VJ <sub>P-</sub>	4.95	8.36	5.62
Dlib <sub>P-</sub>	4.87	8.30	5.55
OD <sub>P-</sub>	4.56	8.16	5.27
GT <sub>P-</sub>	4.43	7.08	5.05
VJ <sub>P</sub>	<b>4.50</b>	<b>7.89</b>	<b>5.16</b>
Dlib <sub>P</sub>	<b>4.42</b>	<b>7.80</b>	<b>5.09</b>
OD <sub>P</sub>	<b>4.36</b>	<b>7.56</b>	<b>4.99</b>
GT <sub>P</sub>	<b>4.36</b>	<b>7.42</b>	<b>4.96</b>

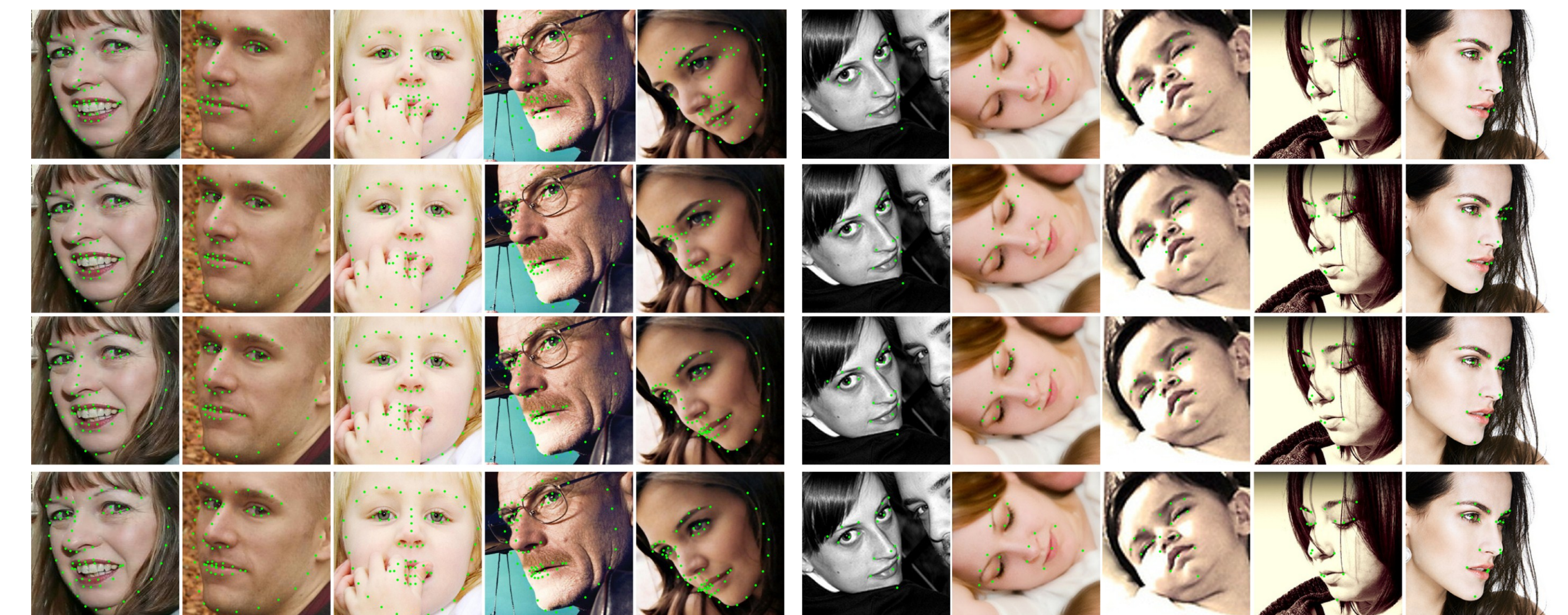
Method	Common Subset	Challenging Subset	Full Set
RCPR [2]	6.18	17.26	8.35
SDM [30]	5.57	15.40	7.52
ESR [5]	5.28	17.00	7.58
CFAN [35]	5.50	16.78	7.69
DeepReg [23]	4.51	13.80	6.31
LBF [21]	4.95	11.98	6.32
CFSS [38]	4.73	9.98	5.76
TCDCN [36]	4.80	8.60	5.54
DDN [34]	-	-	5.59
MDM [25]	4.83	10.14	5.88
Baseline <sub>1</sub>	5.43	8.97	6.12
Baseline <sub>2</sub>	5.03	8.43	5.69
Proposed <sup>-</sup>	4.56	8.16	5.27
Proposed	<b>4.36</b>	<b>7.56</b>	<b>4.99</b>

Comparison with the State-of-the-arts on 300-W

Robustness to Various Initialization

Method	CDM [33]	RCPR	SDM	ERT [16]	LBF	CFSS	CCL [39]	Baseline <sub>1</sub>	Baseline <sub>2</sub>	Proposed <sup>-</sup>	Proposed
NME	5.43	3.73	4.05	4.35	4.25	3.92	2.72	2.99	2.68	2.33	<b>2.17</b>

Comparison with the State-of-the-arts on AFLW



Comparison with the Baselines on 300-W and AFLW