

## Appendix

The appendix is organized as follows:

- In Sec. A1, we provide additional analysis and implementation details of our method.
- In Sec. A2, we provide additional results on WFLW [84], COFW [7], and 300W [66].

### A1. Additional Details

#### A1.1. Additional discussion on Soft-argmax

As mentioned in Sec. 4.1, our work focuses on solving the problem of mismatch with Soft-argmax. The mismatch occurs because the loss is not convex w.r.t. the heatmap’s elements; as illustrated in the example of Eq. (6) and Eq. (7). Our proposed method is based on the structured learning framework of loss-augmented inference [69, 74, 77], where the loss is in the form of a log-sum-of-exponentials, which is convex w.r.t. the heatmap’s elements. In other words, there will be no local minimum with respect to the heatmap, and hence, no mismatch. In Fig. A1, we visualize more examples to show the convergence efficiency comparison between STAR and ours.

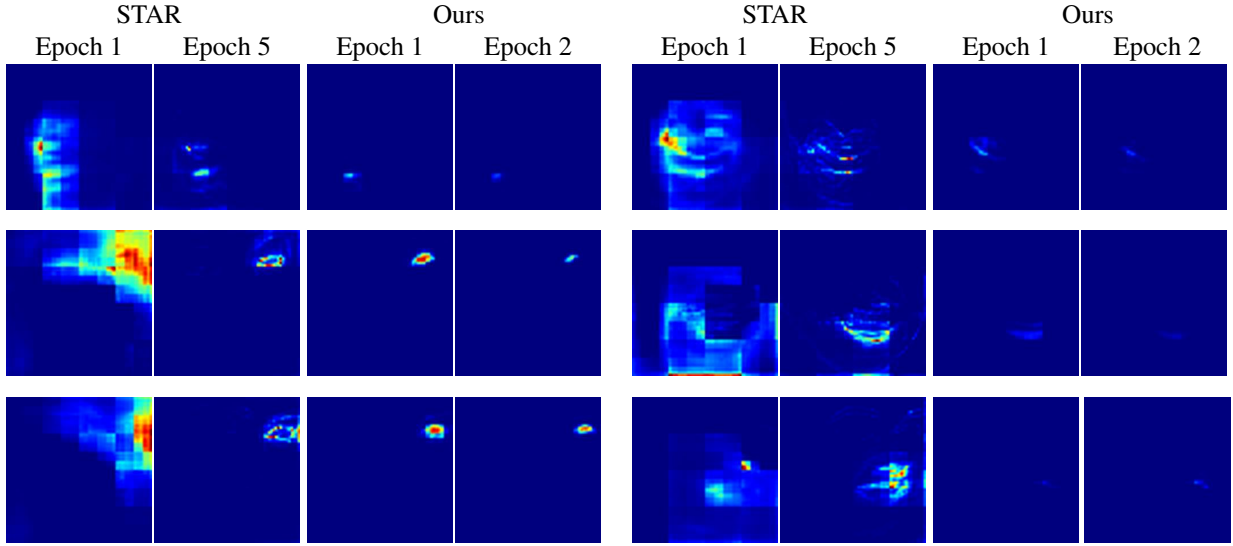


Figure A1. More comparisons of convergence efficiency between STAR’s Soft-argmax and ours following the settings in Fig. 4.

#### A1.2. Additional justification for label smoothing

In Fig. 1, we demonstrate the pipeline of our image-aware label smoothing, which creates a “soft ground-truth”. To verify that this softness resembles semantic ambiguity, in Fig. A2, we ask 5 users to annotate facial images following similar process of WFLW [84]. As expected, the annotations consist of variations. Importantly, we observe a strong similarity between our label smoothing and the distribution of human annotations. This further confirms that our assumption that the variation is along the image edges is correct.

#### A1.3. Implementation Details

**Training details.** The training is conducted on 4 NVIDIA A6000 GPUs with 48GB memory. The training batch size is set to 128. We used the Adam optimizer [39] with a learning rate of 0.001. The model is trained for 200 epochs and takes roughly 10 hours to finish. For the loss function, we follow Huang et al. [35], Zhou et al. [98] and incorporate AAM [35]. AMM is trained on an auxiliary task of predicting edge heatmap ( $\mathcal{L}_{\text{Awing}}$ ) derived from the boundary lines. The total loss function is as follows,

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{Awing}} + \lambda \mathcal{L}_{\text{ours}}, \quad (\text{A14})$$

where  $\mathcal{L}_{\text{Awing}}$  is the AwingLoss [81] for training the AAM [35] and  $\lambda$  controls the weight of our loss  $\mathcal{L}_{\text{ours}}$  from Eq. (13).

**Label smoothing.** Here we provide the implementation details of our label smoothing. Please also refer to Fig. 1. Given  $N$  landmarks  $\mathbf{y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N]$  of an input facial image of size 256px, we follow hand-designed procedure [81, 84]

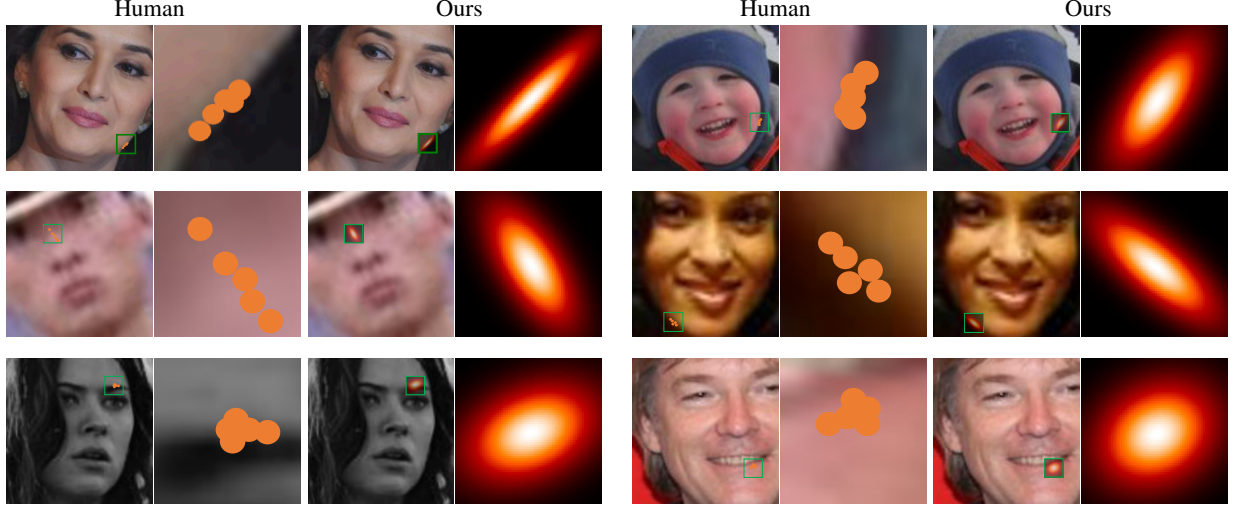


Figure A2. **Comparison between real-life annotation ambiguity and our label smoothing.** The “Human” columns are annotated by 5 users. The “Ours” columns are generated from our label smoothing.

and construct an edge heatmap  $E$  of size 64px. We then apply the `Pytorch` Gaussian blur with a kernel size of 9 and the sharpness filter with a factor of 5 onto  $E$ . The refined edge heatmap is denoted as  $E'$ . For each landmark  $y_n$ , we consider the patch of size  $2k + 1$  around  $y_n$  on  $E'$ , denoted as  $E'_n$ . On the other hand, we construct a Gaussian heatmap  $N_n$  centered around  $y_n$  with a kernel size of 5. After both  $E'_n$  and  $N_n$  are normalized so that their maximum value is 1, we construct a joint heatmap  $M_n = 0.01E'_n + N_n$ . Finally, we build the directional Gaussian heatmap  $G$  for image-aware label smoothing based on the covariance of  $M_n$ , i.e.  $G_n = \mathcal{N}(y_n, \gamma \Sigma_{M_n})$ . The  $\gamma$  is set to 0.001 for COFW but 0.01 for WFLW and 300W.

For each image, we sample 10  $y'_n$  following the distribution  $G$  from our label smoothing. The sampling strategy is lightweight and does not have an observable impact on the overall inference speed.

**Hyperparameters.** We identify the following hyper-parameters which affect the performance of our approach:  $\alpha$  (Eq. 12),  $\epsilon$  (Eq. 11), and  $\lambda$  (Eq. A14). We experiment with different choices, i.e., a grid search, on these hyperparameters to observe their impact on the performance. The results are reported in Fig. A3. It is shown that our method is sensitive to the choices of  $\epsilon$  and  $\lambda$ . Overall, we choose  $\alpha = 1$ ,  $\epsilon = 1$  and  $\lambda = 5000$ .

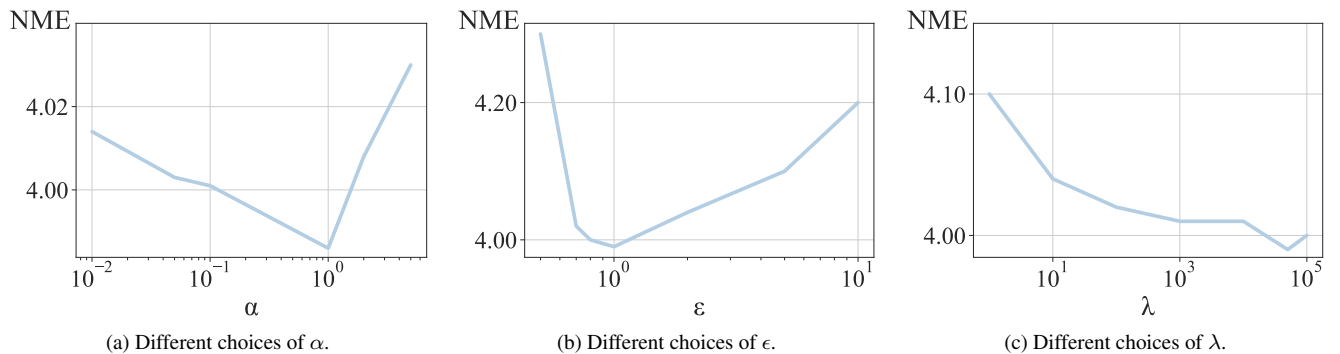


Figure A3. **The NME(↓) of different choices of hyper-parameters on WFLW.**

**Architecture.** Our architecture strictly follows the stacked hourglass network used in STAR and ADNet. Specifically, it has a parameter size of 13.48M and the FLOPs are 17.54G. The throughput is 50 images of 256 pixels by 256 pixels per second.

## A2. Additional Results

In the main paper, we do not use the same evaluation settings across the three datasets presented in Tab. 1 to Tab. 4. This is because not all works report their results on all datasets (WFLW, COFW, and 300W) or all evaluation metrics (NME, FR, and AUC) and settings (inter-ocular or inter-pupil), especially the older methods pre-2020. Additionally, some works do not have open-sourced code or are not reproducible. Note, that our experiment setup follows the state-of-the-art STAR [98], i.e., we compared our method to the same baseline they do. For completeness, we report the comparison of additional evaluation metrics in the following paragraphs.

**Additional ablation studies.** In Tab. A1, we evaluate the robustness of our label smoothing under challenging conditions. We observe a consistent gain on most subsets.

Label smoothing	WFLW-L	WFLW-E	WFLW-I	WFLW-M	WFLW-O	WFLW-B
$\times$	6.68	<b>4.26</b>	3.94	3.93	4.77	<b>4.56</b>
$\checkmark$	<b>6.58</b>	<b>4.26</b>	<b>3.90</b>	<b>3.89</b>	<b>4.74</b>	<u>4.57</u>

Table A1. **Ablation of label smoothing on the subsets of WFLW.** We ablate our label smoothing and report the inter-ocular NME $\downarrow$  on the six subsets of WFLW [8].

**Additional results on WFLW.** In Tab. A2 and Tab. A3, we report the comparison of FR( $\downarrow$ ) and AUC( $\uparrow$ ) to SOTA facial landmark detection methods on the six subsets of WFLW. We achieve competitive results on all subsets. Notably, our method excels in the “largepose” subset. We provide qualitative results of the six subsets in Fig. A4, Fig. A5, Fig. A6, Fig. A7, Fig. A8, and Fig. A9.

Method	Type	WFLW-L	WFLW-E	WFLW-I	WFLW-M	WFLW-O	WFLW-B
Wing [28]	C	22.70	4.78	4.30	7.77	12.50	7.76
LAB [84]	H	28.83	6.37	6.73	7.77	13.72	10.74
DeCaFA [18]	H	21.40	3.73	3.22	6.15	9.26	6.61
AWing [81]	H	13.50	2.23	2.58	2.91	5.98	3.75
ADNet [35]	H	12.72	2.15	2.44	1.94	5.79	3.54
HIH [43]	H	12.88	<b>1.27</b>	2.43	<u>1.45</u>	5.16	3.10
STAR [98]	H	<u>10.77</u>	2.24	<u>1.58</u>	<b>0.98</b>	<u>4.76</u>	<u>2.98</u>
Ours	H	<b>9.23</b>	<u>1.92</u>	<b>0.86</b>	1.46	<b>3.95</b>	<b>2.33</b>

Table A2. **Comparison to SOTA facial landmark detection methods on the subsets of WFLW.** We report the FR $\downarrow$  on the six subsets of WFLW [8], i.e., large pose (WFLW-L), expression (WFLW-E), illumination (WFLWI), make-up (WFLW-M), occlusion (WFLW-O), and blur (WFLW-B). Type “C” and “H” stand for coordinate-regression and heatmap-regression methods.

Method	Type	WFLW-L	WFLW-E	WFLW-I	WFLW-M	WFLW-O	WFLW-B
Wing [28]	C	0.310	0.496	0.541	0.558	0.489	0.491
LAB [84]	H	0.235	0.495	0.543	0.539	0.449	0.463
DeCaFA [18]	H	0.292	0.546	0.579	0.575	0.485	0.494
AWing [81]	H	0.312	0.515	0.578	0.572	0.502	0.512
ADNet [35]	H	0.344	0.523	0.581	0.601	0.530	0.548
HIH [43]	H	0.358	<b>0.601</b>	<b>0.613</b>	0.618	<b>0.539</b>	<b>0.561</b>
STAR [98]	H	<u>0.362</u>	0.584	0.609	<b>0.622</b>	<u>0.538</u>	0.551
Ours	H	<b>0.370</b>	<u>0.587</u>	<u>0.612</u>	<u>0.619</u>	<b>0.539</b>	<u>0.552</u>

Table A3. **Comparison to SOTA facial landmark detection methods on the subsets of WFLW.** We report the AUC $\uparrow$  on the six subsets of WFLW [8], i.e., large pose (WFLW-L), expression (WFLW-E), illumination (WFLWI), make-up (WFLW-M), occlusion (WFLW-O), and blur (WFLW-B). Type “C” and “H” stand for coordinate-regression and heatmap-regression methods.

**Additional results on COFW.** In Tab. A4, we report more quantitative comparisons on COFW. We also provide qualitative results of COFW in Fig. A10.

Method	Inter-Pupil		Inter-Ocular	
	NME ↓	FR ↓	NME ↓	FR ↓
TCDCN [93]	8.05	-	-	-
Wu and Ji [85]	5.93	-	-	-
DAC-CSR [27]	-	-	6.03	4.73
Wing [28]	5.44	3.75	-	-
DCFE [78]	5.27	7.29	-	-
LAB [84]	-	-	3.92	0.39
SDFL [52]	-	-	3.63	<b>0.00</b>
SLPT [86]	4.79	1.18	3.32	<b>0.00</b>
Awing [81]	4.94	0.99	-	-
ADNet [35]	4.68	<u>0.59</u>	-	-
HIH [43]	4.63	<b>0.39</b>	-	-
STAR [98]	<u>4.62</u>	0.79	<u>3.21</u>	<b>0.00</b>
Ours	<b>4.58</b>	0.79	<b>3.15</b>	<b>0.00</b>

Table A4. Additional comparison to SOTA on COFW.

Method	Full	Comm.	Chal.
SDM [88]	7.50	5.57	15.40
CFSS [99]	5.76	4.73	9.98
MDM [76]	5.88	4.83	10.14
RAR [87]	4.94	4.12	8.36
DVLN [83]	4.66	3.94	7.62
HG-HSLE [100]	4.59	3.94	7.24
DCFE [78]	4.55	3.83	7.54
LAB [84]	4.12	3.42	6.98
Wing [28]	4.04	<b>3.27</b>	7.18
Awing [81]	4.31	3.77	6.52
ADNet [35]	4.08	3.51	6.47
STAR [98]	<b>4.03</b>	3.55	6.22
Ours	<b>4.03</b>	3.51	<b>6.09</b>

Table A5. Inter-pupil NME↓ comparisons on 300W and common/challenging subsets.

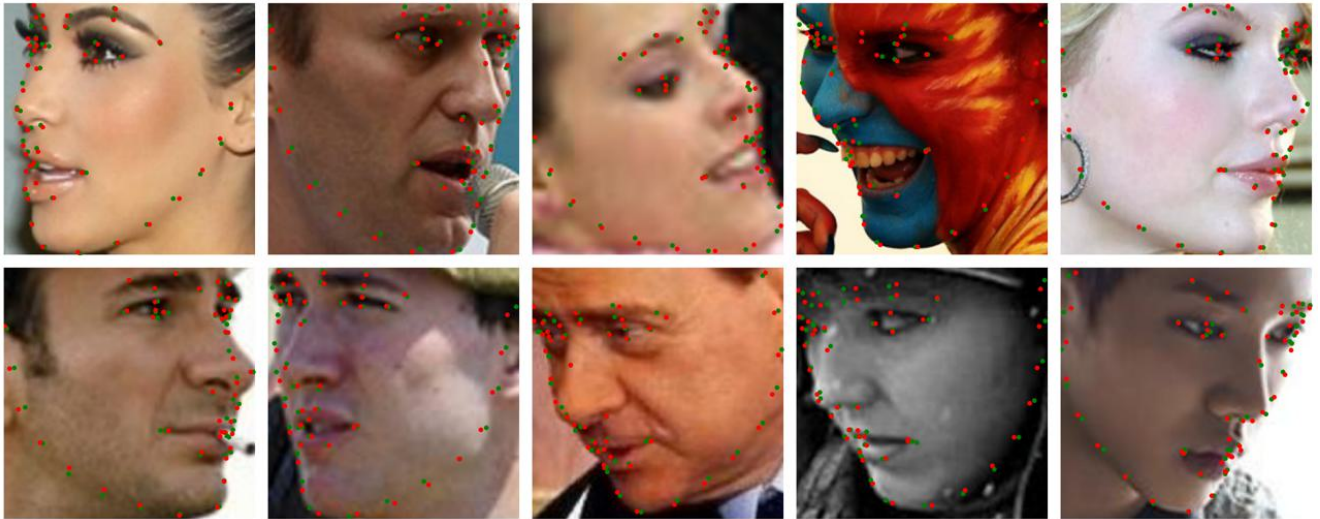


Figure A4. Qualitative results on the “largepose” subset of WFLW. The ground truth landmarks are marked in green while our predictions are marked in red.

**Additional results on 300W.** In Tab. A5, we report more quantitative comparisons on COFW. We provide qualitative results of the two subsets of 300W in Fig. A11 and Fig. A12.



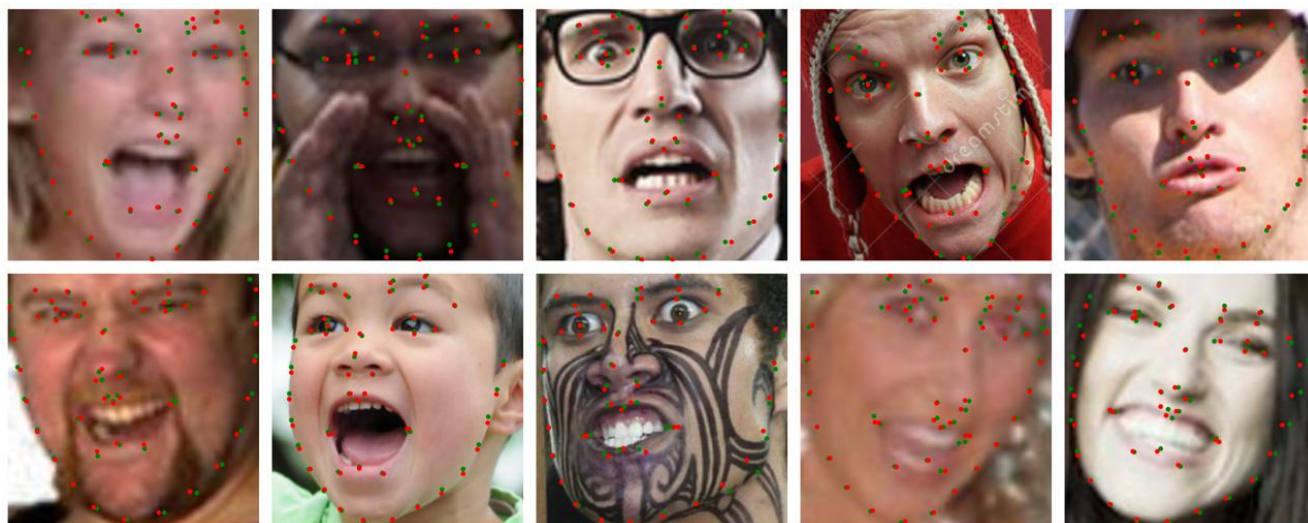


Figure A5. **Qualitative results on the “expression” subset of WFLW.** The ground truth landmarks are marked in **green** while our predictions are marked in **red**.



Figure A6. **Qualitative results on the “illumination” subset of WFLW.** The ground truth landmarks are marked in **green** while our predictions are marked in **red**.



Figure A7. **Qualitative results on the “makeup” subset of WFLW.** The ground truth landmarks are marked in **green** while our predictions are marked in **red**.

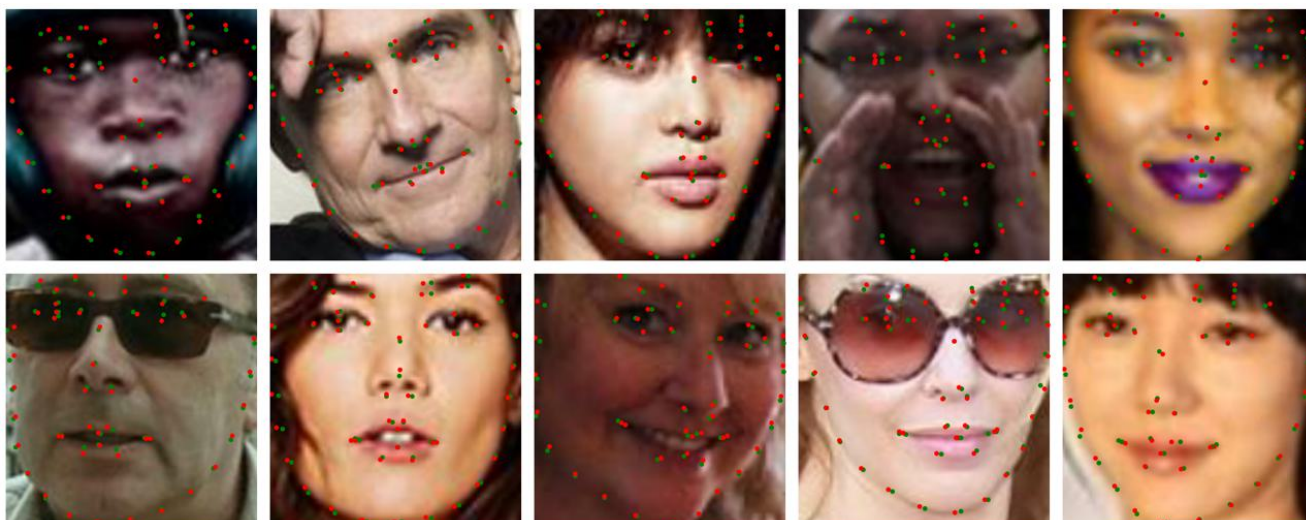


Figure A8. **Qualitative results on the “occlusion” subset of WFLW.** The ground truth landmarks are marked in **green** while our predictions are marked in **red**.



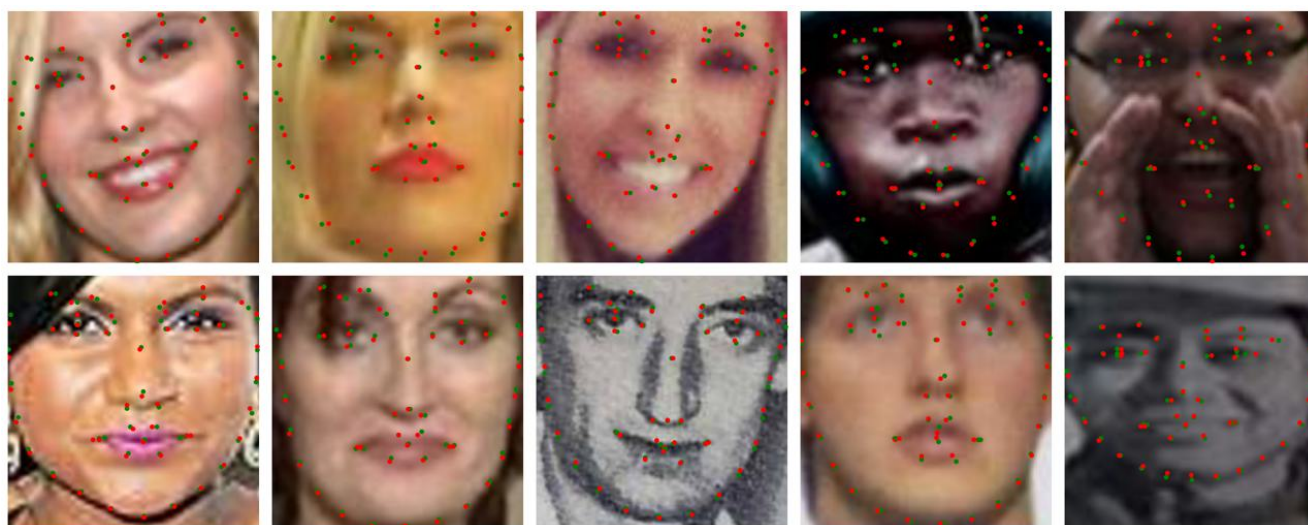


Figure A9. **Qualitative results on the “blur” subset of WFLW.** The ground truth landmarks are marked in **green** while our predictions are marked in **red**.



Figure A10. **Qualitative results on COFW.** The ground truth landmarks are marked in green while our predictions are marked in red.





Figure A11. **Qualitative results on the “common” subset of COFW.** The ground truth landmarks are marked in **green** while our predictions are marked in **red**.



Figure A12. **Qualitative results on the “challenge” subset of 300W.** The ground truth landmarks are marked in **green** while our predictions are marked in **red**.