

Heatmap Regression without Soft-Argmax for Facial Landmark Detection

Chiao-An Yang Raymond A. Yeh
Department of Computer Science, Purdue University

Abstract

Facial landmark detection is an important task in computer vision with numerous applications, such as head pose estimation, expression analysis, face swapping, etc. Heatmap regression-based methods have been widely used to achieve state-of-the-art results in this task. These methods involve computing the argmax over the heatmaps to predict a landmark. Since argmax is not differentiable, these methods use a differentiable approximation, Soft-argmax, to enable end-to-end training on deep-nets. In this work, we revisit this long-standing choice of using Soft-argmax and demonstrate that it is not the only way to achieve strong performance. Instead, we propose an alternative training objective based on the classic structured prediction framework. Empirically, our method achieves state-of-the-art performance on three facial landmark benchmarks (WFLW, COFW, and 300W), converging $2.2\times$ faster during training while maintaining better/competitive accuracy. Our code is available here¹.

1. Introduction

Facial landmark detection, *i.e.*, finding a set of pre-defined landmarks on a given facial image, is an important and classic problem in computer vision. The detected landmarks can aid numerous downstream applications [1, 6, 9, 16, 19, 20, 25, 29, 33, 36, 40, 51, 55, 64, 82, 96], *e.g.*, face recognition, face alignment, face synthesis, head pose estimation, expression analysis, face swapping, *etc.* Research in facial landmark detection has a long history, from earlier works [3, 12–15, 26, 38, 56, 58, 67, 68, 88] to deep learning based methods [52, 71, 86]. Notably, methods based on heatmap regression [17, 18, 22, 29, 35, 37, 41, 43, 48, 50, 57, 63, 80, 81, 84, 92, 98] have found success and gained popularity.

As the name suggests, heatmap regression methods leverage an intermediate representation of heatmaps, one for each landmark. To predict a landmark, these methods compute the

argmax over the coordinates of a given heatmap, *i.e.*, it returns the coordinate that contains the maximum value in the heatmap. However, argmax is not differentiable, hence, Nibali et al. [61] propose Soft-argmax to relax the argmax for end-to-end training. Regression losses, *e.g.*, ℓ_2 -loss, on the predicted coordinate is used to train the model. Next, another challenge in landmark detection is the semantic ambiguity in the annotations [23, 24, 35, 41, 54, 84], meaning that there is some degree of uncertainty in the annotations during the labeling process. To address this issue, several works have studied how to design more robust loss functions [28, 35, 81], *e.g.*, Zhou et al. [98] propose STAR loss, a *self-adaptive* loss to address the semantic uncertainty.

In this work, we revisit the idea of Soft-argmax and find that it may be unnecessary. Instead, we propose an alternative training objective based on the framework of structured prediction that does not require a differentiable coordinate prediction. To address the semantic ambiguity, we propose image-aware label smoothing, which blurs the annotation along the edges of the facial image to simulate label uncertainty. With these two techniques, we achieve a heatmap regression-based method *without* Soft-argmax or a complex loss design.

We evaluate our approach on three established benchmarks, namely, WFLW [84], COFW [7], and 300W [66]. Our approach achieves state-of-the-art performance with a faster training convergence speed while being intuitive and principled. **Our contributions are as follows:**

- We propose to train heatmap regression-based models using a training objective derived from the structured prediction framework and show that the common choice of Soft-argmax may not be necessary.
- We introduce an image-aware label smoothing technique to better capture annotation uncertainty.
- Extensive experiments and ablation studies demonstrate our method’s competitive performance to state-of-the-art while achieving a $2.2\times$ faster training convergence.

2. Related Work

We provide a high-level description of related works on facial landmark detection, structured prediction, and label

¹<https://github.com/ca-joe-yang/regression-without-softarg>

smoothing. Technical details are reviewed in Sec. 3.

Facial landmark detection. The methods for landmark detection can be largely categorized into two branches: coordinates regression and heatmap regression. Coordinate regression methods [5, 28, 45, 46, 76, 94, 97] consider the problem as a regression and directly predict the 2D coordinates. Zhang et al. [94] are the first to leverage additional facial features as auxiliary information. Based on empirical observation, Wing loss [28] is proposed to increase the contribution of predictions with smaller errors. MDM [76] introduces the concept of memory in a coarse-to-fine pipeline. Li et al. [46] employ landmark-to-landmark and landmark-to-memory attention modules to aggregate information to transformer queries. Li et al. [45] experiment with cascaded transformers to explore the structured relationship among facial landmarks.

On the other hand, heatmap regression methods predict landmark heatmaps as intermediate outputs. The coordinates are then acquired by applying the soft approximation [73] of argmax or other functions with similar effects. The main challenge, as opposed to coordinate regression which directly regresses to two numbers, is to be able to predict high-quality landmark heatmaps. Stack hourglass network [60] is initially used in human pose estimation and is later introduced to increase the quality of facial landmark detection [89]. HR-Net [80] maintains high-resolution representations by exchanging the features from different stages. Robinson et al. [65] consider facial landmarks as random variables under Laplace distributions and propose LaplaceKL loss to minimize its KL divergence between prediction and ground truth. Yu and Tao [91] focus on solving the quantization error in heatmap regression. Dapogny et al. [18] consider each stage as an individual alignment task. Auxiliary information, such as facial contours, is also discussed. LAB [84] is the first boundary-aware face alignment algorithm on facial landmark detection.

Besides the landmark heatmap, ADNet [35] generates an edge heatmap and a point heatmap as localization guidance in each hourglass network. It is later refined by STAR [98] to address the annotation ambiguity, applying PCA to the predicted heatmaps to formulate the direction and intensity of the ambiguity. Unlike these heatmap regression works, our method avoids the Soft-argmax approximation by taking inspiration from structured prediction. We also propose an image-aware label smoothing on the *heatmaps annotations* to address semantic ambiguity, whereas Zhou et al. [98] smooths the heatmap predictions. Recently, DISPAL [44] has proposed using the similarity among values of neighboring pixels in addition to the maximum of the heatmap for landmark prediction.

Note, as these works use Soft-argmax , our proposed method can complement their approach. We demonstrate this using STAR [98], the current SOTA with public code.

Structured prediction aims to model relations between the output variables. Seminal works [42, 74, 77] (*e.g.*, Structured Support Vector Machine (SSVM)/ Max-Margin Markov Network) expanded the framework of linear classifiers (*e.g.*, SVM, logistic regression) to consider more complex output structures. Beyond linear models, structure prediction has also been incorporated into deep-nets [4, 10, 11, 21, 31, 32, 49, 53, 70, 75, 90, 95]. A concise tutorial has also been tailored for the computer vision audience [69], with applications in classification and semantic segmentation. In this work, we derive the training objective from the deep structure prediction framework and achieve state-of-the-art performance in the highly competitive task of facial landmark detection.

Label smoothing. Deep-nets are known to overfit and be overconfident in their prediction. One common solution is to train the model using “soft” targets. Instead of a one-hot target with a probability of one in a single class, the target is smoothed with a uniform distribution over all classes [34, 72] or directly penalizing the confidence of a prediction [59, 62]. More advanced label smoothing techniques have also been proposed, *e.g.*, smoothing in a data-dependent way for a classification problem [47]. While our label smoothing is also data-dependent, our smoothing method is inspired by the uncertainty of the annotator; hence, we smooth around the facial boundaries.

3. Background

We briefly review the necessary background and notation to understand this paper.

Landmark detection. Let $X \in \mathbb{R}^{3 \times H \times W}$ to be an input facial image labeled with N landmarks $y = [y_1, y_2, \dots, y_N]$, where each y_n represents the pixel location $(u, v) \in \{0, \dots, H-1\} \times \{0, \dots, W-1\}$ of the landmark, *i.e.*, $y \in \mathbb{R}^{N \times 2}$. The task of landmark detection is to learn a model, parameterized by θ , that predicts landmark \hat{y} given the input face image X .

Heatmap regression. Seminal work by Newell et al. [60] proposes a stacked hourglass architecture and formulates the learning of the model as a heatmap regression task. In heatmap regression methods, given an input, the model outputs N heatmap $\hat{H} = [\hat{H}_1, \dots, \hat{H}_N] \in \mathbb{R}^{N \times H \times W}$ where each heatmap \hat{H}_n represents a score of each pixel being a landmark. To predict a landmark, one will compute the argmax for each of the heatmaps, *i.e.*,

$$\hat{y}_n = \arg \max_{y'_n \in \mathcal{Y}} \hat{H}_n[y'_n], \quad (1)$$

where $\mathcal{Y} \triangleq \{0, \dots, H-1\} \times \{0, \dots, W-1\}$ denotes all the possible pixel locations.

To train the model, Newell et al. [60] create “ground-truth” heatmaps H_n by constructing a 2D Gaussian centered

at each ground-truth landmark \mathbf{y}_n . Training is formulated as minimizing the mean-squared error between the predicted heatmap $\hat{\mathbf{H}}$ and the ground-truth \mathbf{H} :

$$\min_{\theta} \sum_{(\mathbf{X}, \mathbf{y}) \in \mathcal{D}} \|\hat{\mathbf{H}}(\mathbf{X}; \theta) - \mathbf{H}(\mathbf{y})\|_2^2, \quad (2)$$

where \mathcal{D} denotes the training dataset.

While this formulation successfully learns a model, it is unclear whether using mean-square error on heatmaps is a suitable loss. Importantly, the evaluation metric of landmark detection compares between the predicted and ground-truth *landmark* and *not* heatmaps. One idea is to use a differentiable argmax such that the models can be trained using a loss function on the landmarks, which we discuss next.

Soft-argmax. To enable training using the gradient method on losses between the predicted and ground-truth landmarks, Nibali et al. [61] propose a soft-argmax operation. Given a heatmap \mathbf{H}_n , the Soft-argmax normalizes it into a probability distribution and computes the expectation over the coordinates. Formally, $\tilde{\mathbf{y}}_n =$

$$\text{Soft-argmax}(\hat{\mathbf{H}}_n) \triangleq \sum_{\mathbf{y}'_n} \mathbf{y}'_n \cdot \text{Softmax}(\hat{\mathbf{H}}_n)[\mathbf{y}'_n], \quad (3)$$

where Softmax normalizes over all the spatial dimensions to form a probability distribution. With Soft-argmax defined, the model can be trained by minimizing a loss, e.g., ℓ_2 -loss as follows:

$$\min_{\theta} \sum_{(\mathbf{X}, \mathbf{y}) \in \mathcal{D}} \sum_n \|\tilde{\mathbf{y}}_n - \mathbf{y}_n\|_2^2. \quad (4)$$

While soft-argmax has been widely used [43, 98], we re-examine whether such relaxation is a good design choice.

Deep structured learning. The design of training objective over an argmax in Eq. (1) is well studied in structured prediction. A typical choice for linear models is to use the structural hinge loss [74], leading to Structured-SVMs [77].

Chen et al. [11], Schwing and Urtasun [69] further generalize the idea and present deep structured learning with a training objective:

$$\min_{\theta} \frac{C}{2} \|\theta\|_2^2 + \sum_{(\mathbf{X}, \mathbf{y}) \in \mathcal{D}} \underbrace{\max_{\hat{\mathbf{y}}} (\Delta(\mathbf{y}, \hat{\mathbf{y}}) + F(\hat{\mathbf{y}}, \mathbf{X}; \theta)) - F(\mathbf{y}, \mathbf{X}; \theta)}_{\text{Loss-augmented inference}}, \quad (5)$$

where Δ denotes the margin term (*a.k.a.* task-loss), and F denotes a score function represented by a deep-net with parameters θ , and $C \in \mathbb{R}$ is a hyper-parameter controlling the regularization term. Note, in the case where the score is a linear model, *i.e.*, a “one layer deep-net” with the choice $F = \theta^\top \text{vec}(\mathbf{X})$, then Eq. (5) recovers the objective of a Structured-SVM.

In this work, we demonstrate that this framework provides a principled alternative to the existing training methods for heatmap regression in Eq. (2) or soft-argmax in Eq. (4) commonly used by prior works [35, 43, 98].

4. Approach

We first analyze the shortcomings of Soft-argmax. We then describe how to formulate landmark detection with the deep structured learning framework along with a label smoothing technique to handle noisy annotations.

4.1. Limitation of Soft-argmax

As reviewed in Sec. 3, Soft-argmax approximates the argmax with the expected value of the argument over a distribution form by normalizing a heatmap. Below, we show why this approximation may not be desirable. For readability, we illustrate the example using a “1-D heatmap”. Let the ground-truth landmark $\mathbf{y} = 2$, and consider two heatmaps in Eq. (6) and Eq. (7).

Soft-argmax $\tilde{\mathbf{y}}^{(1)} = 2$

$\mathbf{y}'^{(1)}$	0	1	2	3	4
Softmax($\hat{\mathbf{H}}^{(1)}$)	0.0	0.0	1.0	0.0	0.0

(6)

Unimodal heatmap

Soft-argmax $\tilde{\mathbf{y}}^{(2)} = 2$

$\mathbf{y}'^{(2)}$	0	1	2	3	4
Softmax($\hat{\mathbf{H}}^{(2)}$)	0.4	0.1	0.0	0.1	0.4

Bimodal heatmap

(7)

Observation: Both the unimodal and bimodal heatmap results in the same Soft-argmax output $\tilde{\mathbf{y}}^{(1)} = \tilde{\mathbf{y}}^{(2)} = 2$. In other words, the ℓ_2 -loss $\|\tilde{\mathbf{y}}^{(2)} - \mathbf{y}\| = 0$ which means the training has converged. However, the **bimodal heatmap makes an incorrect argmax prediction** as $\{0, 4\} = \arg \max_{\mathbf{y}'} \hat{\mathbf{H}}^{(2)}[\mathbf{y}'] \neq 2$. More cases with a mismatch can be easily found, e.g., when there are even more peaks. We believe this mismatch may lead to optimization difficulties.

We empirically validate the hypothesis and provide empirical study in our experiments (see Sec. 5.1). With this observation in mind, we rethink whether Soft-argmax is necessary to train an effective model. More analysis is provided in Appendix A1.1. We now discuss our proposed training objective based on structured prediction.

4.2. Deep Structured Landmark Detection

Problem formulation. As reviewed in Sec. 3, given an image \mathbf{X} , the model aims to predict the ground truth landmarks \mathbf{y} consist N individual landmarks \mathbf{y}_n , each corresponding to an (u, v) 2D landmark coordinates, e.g., the top of the lips.

We propose to formulate the training of a landmark detection model as deep structured learning in Eq. (5). The

key question is how to design the score function $F(\mathbf{y}, \mathbf{X}; \theta)$. First, we choose a score function that is decomposed into a sum of score functions F_n for each landmark, *i.e.*,

$$F(\mathbf{y}, \mathbf{X}; \theta) \triangleq \sum_n F_n(\mathbf{y}_n, \mathbf{X}; \theta), \quad (8)$$

where $F_n(\mathbf{y}_n, \mathbf{X}; \theta)$ corresponds to a channel of the stacked hourglass architecture’s output. More precisely, the score $F_n(\mathbf{y}_n, \mathbf{X}; \theta) \triangleq \hat{H}_n[\mathbf{y}_n]$.

Substituting the function of F into Eq. (5) and simplify, we arrive at the training objective:

$$\min_{\theta} \frac{C}{2} \|\theta\|_2^2 + \sum_{(\mathbf{X}, \mathbf{y}) \in \mathcal{D}} \sum_{n=1}^N \max_{\hat{\mathbf{y}}_n} (\Delta(\mathbf{y}_n, \hat{\mathbf{y}}_n) + F_n(\hat{\mathbf{y}}_n, \mathbf{X}; \theta)) - F_n(\mathbf{y}_n, \mathbf{X}; \theta). \quad (9)$$

Next, the objective can be relaxed using the fact that

$$\epsilon \ln \sum_i \exp \frac{x_i}{\epsilon} \xrightarrow{\epsilon \rightarrow 0} \max_i x_i \quad (10)$$

to obtain the final training objective

$$\min_{\theta} \frac{C}{2} \|\theta\|_2^2 + \sum_{(\mathbf{X}, \mathbf{y}) \in \mathcal{D}} \sum_n \mathcal{L}(\mathbf{X}, \mathbf{y}_n, \theta), \quad (11)$$

where we subsume all terms into a loss $\mathcal{L}(\mathbf{X}, \mathbf{y}_n, \theta)$ as

$$\epsilon \ln \left(\sum_{\hat{\mathbf{y}}_n} \exp \frac{\Delta(\mathbf{y}_n, \hat{\mathbf{y}}_n) + F_n(\hat{\mathbf{y}}_n, \mathbf{X}, \theta)}{\epsilon} \right) - F_n(\mathbf{y}_n, \mathbf{X}, \theta).$$

For an effective margin Δ , it needs to capture some notion of “similarity” for the task, *i.e.*, how close is the prediction $\hat{\mathbf{y}}_n$ to ground-truth \mathbf{y}_n and satisfy $\Delta(\mathbf{y}, \mathbf{y}) = 0$. For simplicity, we choose the margin term to be a standard smooth- ℓ_1 distance [30], *i.e.*,

$$\Delta(\mathbf{y}_n, \hat{\mathbf{y}}_n) = \alpha \cdot \begin{cases} \frac{0.5}{s} \|\hat{\mathbf{y}}_n - \mathbf{y}_n\|_2^2, & \text{if } \|\hat{\mathbf{y}}_n - \mathbf{y}_n\|_1 < s, \\ \|\hat{\mathbf{y}}_n - \mathbf{y}_n\|_1 - 0.5s, & \text{otherwise.} \end{cases} \quad (12)$$

This is a widely used distance function in previous work [35, 98]. Here, the threshold s is set as 0.01 and α is a weighting factor controlling the significance of the margin term.

Finally, any gradient-based optimization methods, *e.g.*, Adam [39], can be used to optimize the training objective in Eq. (11). We note that as $\hat{\mathbf{y}}_n$ is not a function of θ , no gradients are computed through the margin term Δ .

Image-aware label smoothing. The annotation for landmark detection is noisy due to semantic ambiguity [41, 54, 84]. To handle this uncertainty, we propose to use label-smoothing. Instead of having a single “one-hot” annotation, we construct a smooth distribution \mathbf{G} over it, *i.e.*, the training objective in Eq. (11) is smoothed into

$$\min_{\theta} \frac{C}{2} \|\theta\|_2^2 + \sum_{(\mathbf{X}, \mathbf{y}) \in \mathcal{D}} \sum_n \mathbb{E}_{\mathbf{y}'_n \sim \mathbf{G}(\mathbf{y}_n)} [\mathcal{L}(\mathbf{X}, \mathbf{y}'_n, \theta)], \quad (13)$$

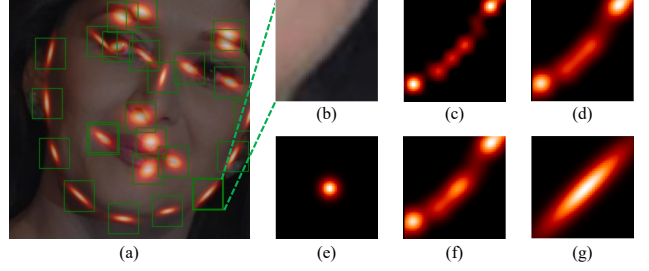


Figure 1. **Our label smoothing.** We visualize the pipeline of our implementation on label smoothing. (a) is the overall label smoothing result. (b-g) are intermediate products cropped around the selected ground truth landmark \mathbf{y} . (b) is the input image patch. (c) is the unprocessed pseudo edge heatmap. (d) is the processed edge heatmap. (e) is a Gaussian heatmap centered around \mathbf{y} . (f) is the weighted summation of (d) and (e). (g) is the Gaussian heatmap based on the covariance of (f).

where we approximate the expected value with Monte Carlo samples.

As in prior works of semantic ambiguity, the insight into constructing this smooth distribution \mathbf{G} comes from the observation that the annotation *varies along* the local edges around the landmark (see Fig. 1 (a), (b)). To construct \mathbf{G} , we leverage the pseudo edge heatmap (shown in Fig. 1 (c)) generated from the landmark annotations. To create such a pseudo edge heatmap, we follow a hand-designed procedure in prior works [81, 84]. This involves interpolating landmark annotations to create boundary lines followed by a distance transform which weights each pixel by its distance to the boundary lines.

Next, we notice the edge heatmap is often disjointed; hence, we apply blurring and sharpening transformations to enhance the quality (Fig. 1 (d)). To ensure that the ground-truth location is the maximum location, we add a normal Gaussian distribution (Fig. 1 (e)) centered at each landmark, resulting in Fig. 1 (f). Finally, we extract a Gaussian distribution centered at \mathbf{y}_n with a covariance estimated from (f) to be used as our smoothed distribution \mathbf{G} . Please note that the pseudo edge heatmap is generated from ground truth landmark annotations and was also used by the compared baselines [35, 98], hence, no additional data was used.

5. Experiments

We start with a simplified model to analyze the differences between our training objective versus ℓ_2 -loss + Soft-argmax. We then conduct experiments on three standard landmark detection benchmarks, followed by ablation studies. Finally, we show that our loss achieves significantly faster convergence speed.

5.1. Analysis on a simplified model

We consider a simplified model with a single landmark on a 1D heatmap. We choose the dataset to have a single

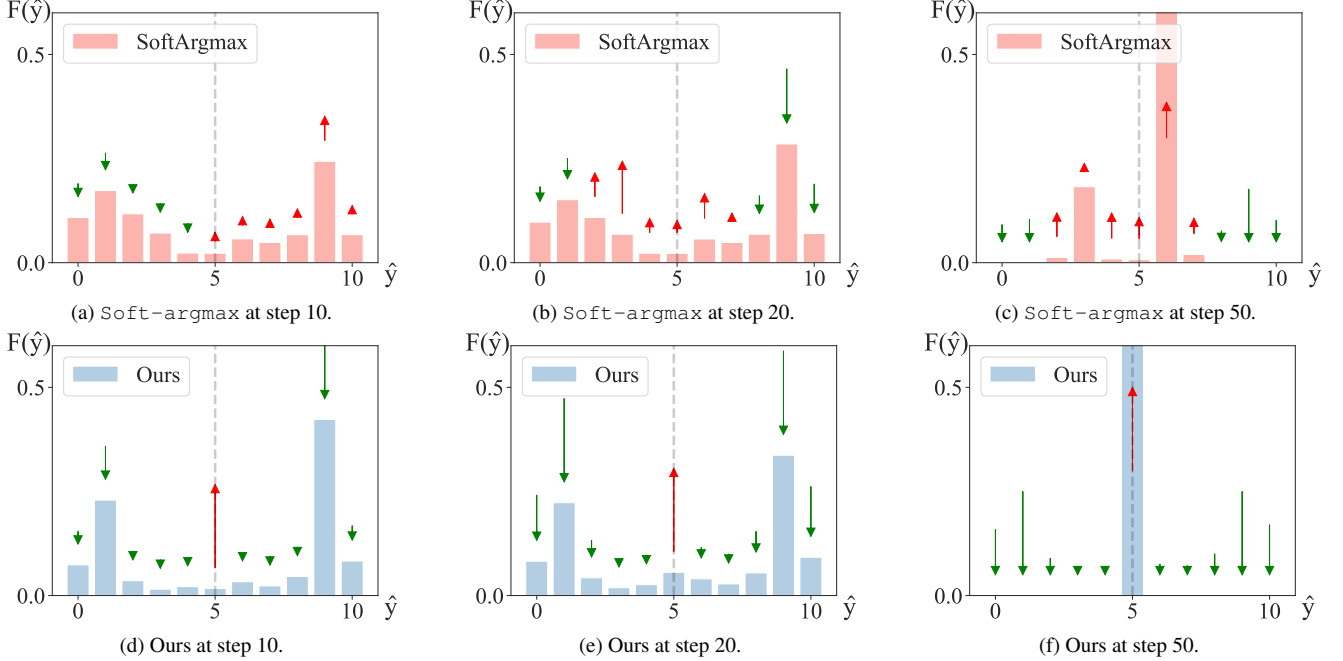


Figure 2. **Comparison between Soft-argmax and ours.** (a-c) visualize each F trained with Eq. (4), i.e. Soft-argmax and ℓ_2 -loss at step 10, 20, and 50. (d-f) visualize each F trained with our loss at steps 10, 20, and 50. We use \uparrow to denote where F is increasing, i.e., positive updates, and \downarrow to denote where F is decreasing, i.e., negative updates.

training sample at $y = 5$, in this case, we can remove the dependencies on the input X and directly model the heatmap with the model parameters. That is, the heatmap $\hat{H}[k] \triangleq \theta_k$ is parameterized by $\theta \in \mathbb{R}^{11}$. In Fig. 2, we visualize how gradient descent updates the heatmap, i.e., $\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta \nabla_{\theta} \mathcal{L}$, when optimizing with the objective ℓ_2 -loss + Soft-argmax (Eq. (4)), or our proposed method (Eq. (11)). We initialize $\theta^{(0)}$ to be a bimodal distribution with $\theta_9^{(0)}$ and $\theta_1^{(0)}$ each having the highest value and second highest values. For our method, we choose the margin term Δ to be the ℓ_2 -loss.

In Fig. 2 (a-c), we show how θ changes when training with ℓ_2 -loss + Soft-argmax (Eq. (4)) with arrows showing the update. As can be observed, θ_5 is being updated in the correct direction, however, its update magnitude is fairly small. Furthermore, the updates for the other parameters are not consistent. Remember, as $y = 5$, the heatmap should have θ_5 being the largest element.

In Fig. 2 (d-f), we show how θ changes when training with our proposed method (Eq. (11)). We observe that the updates are straightforward. The parameter θ_5 consistently receives positive updates, while all other parameters receive negative updates. After a few steps, the argmax is already located $\hat{H}[5]$, with a significant gap between its maximum and second maximum. That is, we observe faster convergence of our approach. Overall, we observe that ℓ_2 -loss + Soft-argmax gradually “moves” the peaks towards the target $y = 5$. In contrast, our proposed approach (Eq. (11))

Method	Type	NME \downarrow	FR \downarrow	AUC \uparrow
Wing [28]	C	4.99	6.00	0.550
SLPT [86]	C	4.14	2.76	0.595
RePFormer [46]	C	4.11	-	-
DTLD [45]	C	4.08	2.76	-
DeCaFa [18]	H	4.62	4.84	0.563
HRNet [80]	H	4.60	4.64	-
LUVLi [41]	H	4.37	3.12	0.577
Awing [81]	H	4.36	2.84	0.572
PIPNNet [37]	H	4.31	-	-
DAG [48]	H	4.21	3.04	0.589
ADNet [35]	H	4.14	2.72	0.602
HIH [43]	H	4.08	2.60	<u>0.605</u>
KeyPosS [2]	H	<u>4.00</u>	-	-
STAR [98]	H	4.02	2.32	0.605
Ours	H	3.99	1.84	0.606

Table 1. **Comparison to SOTA facial landmark detection methods on the full set of WFLW.** We report the inter-ocular NME \downarrow , FR $_{10\%}\downarrow$, and AUC $_{10\%}\uparrow$ on WFLW [8]. “C” and “H” correspond to coordinate/heatmap-regression, respectively.

directly increases the peak at the target $y = 5$ while decreasing the others.

5.2. Landmark detection

Experiment setup. We conduct experiments on commonly used landmark detection benchmarks:

Method	Type	WFLW-L	WFLW-E	WFLW-I	WFLW-M	WFLW-O	WFLW-B
Wing [28]	C	8.75	5.36	4.93	5.41	6.37	5.81
RePFormer [46]	C	7.25	<u>4.22</u>	4.04	3.91	5.11	4.76
DeCaFa [18]	H	8.11	4.65	4.41	4.63	5.74	5.38
AWing [81]	H	7.38	4.58	4.32	4.27	5.19	4.96
ADNet [35]	H	6.96	4.38	4.09	4.05	5.06	4.79
HIH [43]	H	6.87	4.06	4.34	<u>3.85</u>	4.85	4.66
STAR [98]	H	<u>6.76</u>	4.27	<u>3.97</u>	3.83	<u>4.80</u>	<u>4.58</u>
Ours	H	6.58	4.26	3.90	3.89	4.74	4.57

Table 2. **Comparison to SOTA facial landmark detection methods on the subsets of WFLW.** We report the inter-ocular NME \downarrow on the six subsets of WFLW [8], *i.e.*, large pose (WFLW-L), expression (WFLW-E), illumination (WFLW-I), 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	NME \downarrow	FR \downarrow	AUC \uparrow
DAC-CSR [27]	C	6.03	4.73	-
Wu and Ji [85]	C	5.93	-	-
Wing [28]	C	5.44	3.75	-
DCFE [78]	C	5.27	7.29	0.359
SLPT [86]	C	4.79	1.18	-
MHHN [79]	H	4.95	1.78	-
Awing [81]	H	4.94	0.99	0.488
ADNet [35]	H	4.68	<u>0.59</u>	0.532
HIH [43]	H	4.63	0.39	-
STAR [98]	H	<u>4.62</u>	0.79	0.540
Ours	H	4.58	0.79	0.544

Table 3. **Comparison to SOTA on COFW.** “C” and “H” stand for coordinate-regression and heatmap-regression methods.

- WFLW [84] contains 7,500 training and 2,500 test images, each labeled with 98 facial landmarks. WFLW contains six subsets, including “large pose”, “expression”, “illumination”, “make-up”, “occlusion”, and “blur”. These subsets identify the more challenging conditions for facial landmark detection.
- COFW [7] contains 1,345 training and 507 test images, each labeled with 68 facial landmarks. 300W contains 3,148 and 689 test images, each labeled with 29 facial landmarks.
- 300W [66] test images are split into two: a common subset (Comm.) with 554 images and a challenging (Chal.) subset with 135 images under more diverse illumination conditions and variations of expressions and poses.

Evaluation metric. As in prior works, we report on three metrics: (a) *Normalized Mean Error (NME \downarrow)*, (b) *Failure Rate (FR \downarrow)*, and (c) *Area Under Curve (AUC \uparrow)*. NME is defined as the average ℓ_2 -loss over the landmarks: $NME(y_n, \hat{y}_n) \triangleq \frac{1}{N \cdot d} \sum_{n=1}^N \|y_n - \hat{y}_n\|_2$, where d is the distance used for normalizing the error.

The FR is the ratio of “failed” cases, where a case is con-

Method	Type	Full	Comm.	Chal.
SLPT [86]	C	3.17	2.75	4.90
RePFormer [46]	C	3.01	-	-
DTLD [45]	C	2.96	2.59	4.50
DeCaFa [18]	H	3.39	2.93	5.26
HRNet [80]	H	3.32	2.87	5.15
HIH [43]	H	3.09	2.65	4.89
Awing [81]	H	3.07	2.72	4.52
ADNet [35]	H	2.93	<u>2.53</u>	4.58
KeyPosS [2]	H	3.34	-	-
STAR [98]	H	2.87	2.52	<u>4.32</u>
Ours	H	2.87	<u>2.53</u>	4.27

Table 4. **Inter-ocular NME \downarrow comparisons on 300W and common/challenging subsets.**

sidered to fail if its NME exceeds a given threshold. The AUC measures the area under the Cumulative Error Distribution (CED) curve from zero to a given NME threshold. Following prior works [35, 98], we use inter-ocular distance as d for both WFLW and 300W and inter-pupil distance for COFW. The FR and AUC threshold is set to 10% for WFLW and COFW, and 5% for 300W.

Implementation details. As we are innovating on the loss function, we strictly follow the data preparation and the deep-net architecture from STAR by Zhou et al. [98]. During training, each input image first has the face region cropped out and resized to 256px \times 256px. The image is then processed through the following random data augmentations, including rotation, scaling, brightness, blur, cutout, and horizontal flip.

For the deep-net architecture, we use the same four-stacked hourglass architecture [98] to predict the heatmap. We replaced their Soft-argmax layer along with ReLU which they used to normalize the heatmap. Note that STAR [98] uses the Anisotropic Attention Module (AAM), an attention mask to guide landmark prediction. The AAM is trained with an auxiliary task of predicting edge heatmaps

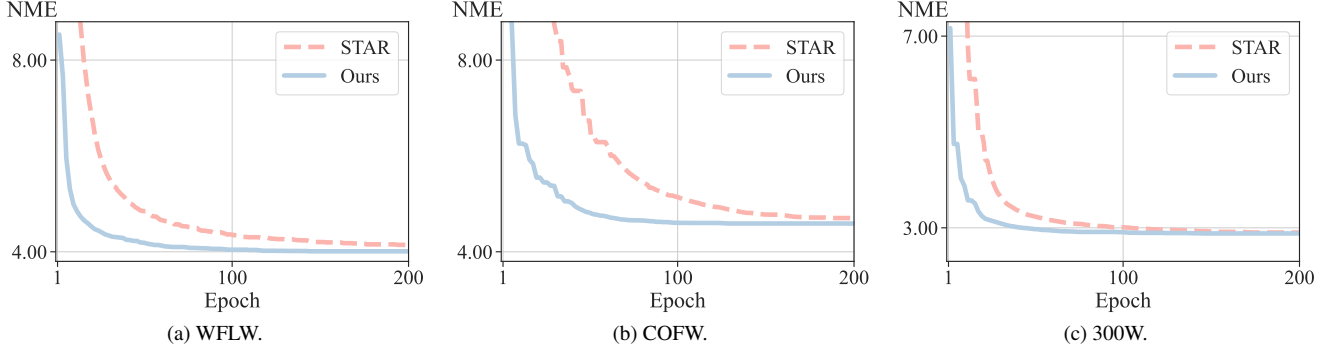


Figure 3. **Comparison of NME curve throughout training between STAR and ours.** We plot the curve of inter-ocular NME \downarrow over training epochs for both methods. The performance improves and converges much faster using our method than using STAR [98].

derived from the boundary lines. We also included such auxiliary tasks in our training. Additional implementation details and hyperparameters are documented in the Appendix Sec. A1.

Quantitative results. For baselines, we compare to the SOTA facial landmark detection methods with available code or reproducibility. In Tab. 1, we report the quantitative comparison on the full test set of WFLW. The best result is **bolded** and the second best is underlined. Our method achieves state-of-the-art performance under all evaluation metrics. More notably, we outperform STAR [98] by 0.02 on NME and 0.32 on FR, which are comparable to the gains achieved by prior works. Note, this result is significant. We trained our model three times and the NME’s standard deviation is 0.003.

In Tab. 2, we report the results on the six more challenging subsets of WFLW. We achieve state-of-the-art on 4 subsets while being competitive on the remaining 2 subsets. Notably, our method outperforms STAR [98] the most under the “large-pose” category (WFLW-L), where the input facial images have larger motion and the poses are far from frontal. We achieve an NME improve of 0.18.

In Tab. 3, we report the comparison between baselines and ours on COFW. Our method achieves state-of-the-art performance in terms of NME and AUC on COFW.

In Tab. 4, we report the performance under the subsets of 300W. Again, our method remains competitive to SOTA under all categories. Additional results are in Sec. A1.

Training convergence. In Fig. 3, we visualize the learning efficiency of our method in comparison to STAR [98]. We plot the training curve, *i.e.*, NME versus training epochs, on WFLW, COFW, and 300W. Consistently across these experiments, we observe a much faster convergence speed with our method. Roughly speaking, our loss converges $2.2\times$ faster than STAR. For example, it takes STAR 44 epochs to achieve an NME of 4.50 on WFLW and ours to achieve the same performance in 20 epochs.

In Fig. 4, we further demonstrate the convergence lim-

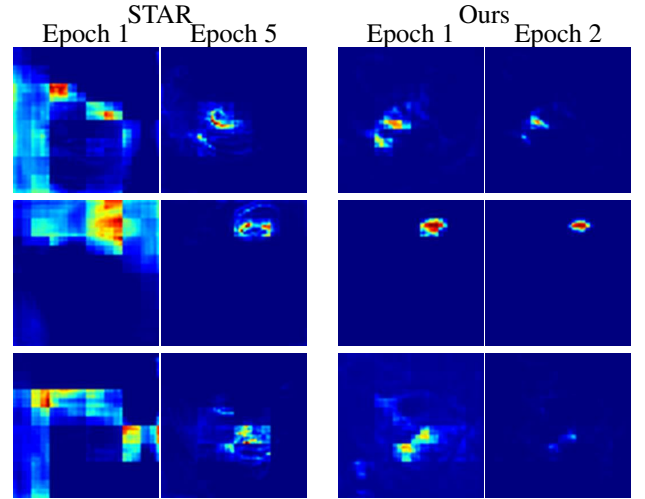


Figure 4. **Comparison of convergence efficiency between STAR’s Soft-argmax and ours.** The first two columns visualize the \hat{H} trained with STAR [98]’s Soft-argmax at epochs 1 and 5. The last two columns visualize the \hat{H} trained with our loss at epochs 1 and 2. Each row corresponds to a different sample.

itations of Soft-Argmax, we visualize the predictions \hat{H} trained with STAR’s Soft-Argmax versus ours on WFLW. We observe the same behavior as in Fig. 2. STAR’s Soft-Argmax still suffers from bi-modal prediction at epoch $t = 5$, while ours quickly finds the correct peak in 2 epochs. More comparisons are provided in the Appendix A1.1.

Qualitative results. In Fig. 5, we provide a qualitative comparison between STAR and Ours. We observe comparable results in the predicted key points. In particular, we observe slight benefits along the facial contours, *e.g.*, jaw lines, and in more challenging cases involving occlusions.

5.3. Ablation studies

Choices for the margin term Δ . In Tab. 5, we report the performance of other common distance functions in facial

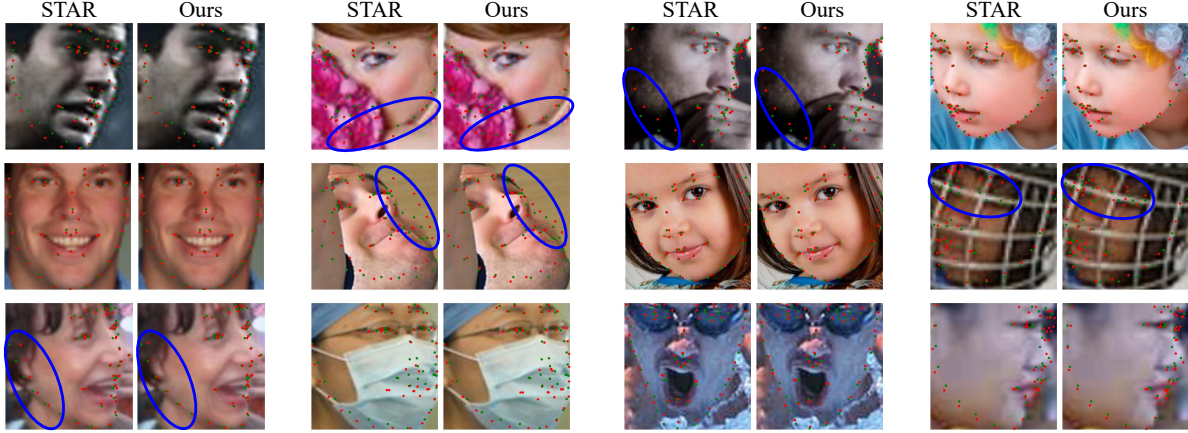


Figure 5. **Qualitative comparison on WFLW between STAR and ours.** The ground truth landmarks are marked in green while the predictions are marked in red. The regions highlighted in blue circles emphasize where our method outperforms STAR [98].

Δ	NME↓	FR↓	AUC↑
None	4.01	2.16	0.605
ℓ_2	4.00	2.00	0.606
ℓ_1	4.00	2.04	0.606
smooth- ℓ_1	3.99	1.84	0.606

Table 5. **Ablation on the choice of margin Δ .** Smooth- ℓ_1 achieves the best performance.

landmark detection as the margin term Δ . We observe that the choices of Δ mildly impact the final performance. We note that by removing the margin term, Eq. (11) can be reduced to cross-entropy loss with a temperature term. Overall, our choice of Δ of smooth- ℓ_1 [30] has the best performance.

Structure prediction and label smoothing. In Tab. 6, we report additional ablation studies on our loss. We consider different configurations for label smoothing, *i.e.*, no label smoothing, image-unaware smoothing, and image-aware smoothing. We also consider different configurations for loss design, *i.e.*, STAR [98] and ours based solely on structured prediction in Eq. (11). Specifically, the rows represent the following configurations: (a) no contribution; (b) only structured prediction; (c) only image-unaware label smoothing; (d) only image-aware label smoothing; (e) both structured prediction and image-aware label smoothing.

Without label smoothing, our loss achieves comparable performance to STAR [98] on WFLW. We can also observe improvement when applying label smoothing on STAR. We also observe a consistent gain in deploying label smoothing under all evaluation metrics and being image-aware further improves performance. Finally, the best performance is achieved by incorporating all of our contributions.

Importantly, by incorporating image-aware label smoothing, we outperform without label smoothing by 0.02 on NME, 0.23 on FR, and 0.003 on AUC. Validating its effec-

Label Smoothing	Loss	NME↓	FR↓	AUC↑
\times	STAR	4.02	2.32	0.605
\times	Eq. 11	4.02	2.23	0.604
image-unaware	Eq. 11	3.99	<u>2.12</u>	0.606
image-aware	STAR	4.00	2.12	0.602
image-aware	Eq. 11	3.99	1.84	0.606

Table 6. **Ablation on structure prediction and label smoothing.**

tiveness.

Limitations: While the proposed label smoothing works well, its design is hand-crafted and more of a heuristic based on our observation. In the future, we aim to explore data-driven smoothing approaches and study how to better quantify the uncertainty.

6. Conclusion

In this work, we re-examine and analyze the widely used Soft-argmax in recent heatmap regression methods for facial landmark detection. Instead, we propose to formulate the training objective based on deep structured prediction, for which we can more effectively train the model without Soft-argmax. To further address semantic ambiguity, we propose an image-aware label-smoothing technique. With these two components, we arrived at a model with a principled training objective. Empirically, our method obtained state-of-the-art performance on three common facial landmark detection datasets with $2.2\times$ faster training convergence. We hope this work will inspire the revisiting of popular design choices in other computer vision tasks through the lens of structured prediction, which may lead to cleaner and more intuitive designs.

Acknowledgments. RAY is thankful for the ECE544 course at UIUC, taught by Alex Schwing, which provided the foundational concepts for this work.

References

- [1] Jianmin Bao, Dong Chen, Fang Wen, Houqiang Li, and Gang Hua. Towards open-set identity preserving face synthesis. In *Proc. CVPR*, 2018. 1
- [2] Xu Bao, Zhi-Qi Cheng, Jun-Yan He, Wangmeng Xiang, Chenyang Li, Jingdong Sun, Hanbing Liu, Wei Liu, Bin Luo, Yifeng Geng, et al. KeyPosS: Plug-and-play facial landmark detection through gps-inspired true-range multilateration. In *Proc. ACM MM*, 2023. 5, 6
- [3] Volker Blanz and Thomas Vetter. Face recognition based on fitting a 3d morphable model. *IEEE TPAMI*, 2003. 1
- [4] John Bridle. Training stochastic model recognition algorithms as networks can lead to maximum mutual information estimation of parameters. In *Proc. NeurIPS*, 1989. 2
- [5] Bjorn Browatzki and Christian Wallraven. 3fabrec: Fast few-shot face alignment by reconstruction. In *Proc. CVPR*, 2020. 2
- [6] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem?(and a dataset of 230,000 3d facial landmarks). In *Proc. ICCV*, 2017. 1
- [7] Xavier Burgos-Artizzu, Pietro Perona, and Piotr Dollar. Caltech occluded faces in the wild (COFW), 2022. 1, 6, 12
- [8] Xavier P Burgos-Artizzu, Pietro Perona, and Piotr Dollár. Robust face landmark estimation under occlusion. In *Proc. ICCV*, 2013. 5, 6, 14
- [9] Xudong Cao, Yichen Wei, Fang Wen, and Jian Sun. Face alignment by explicit shape regression. *IJCV*, 2014. 1
- [10] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Semantic image segmentation with deep convolutional nets and fully connected CRFs. In *Proc. ICLR*, 2015. 2
- [11] Liang-Chieh Chen, Alexander Schwing, Alan Yuille, and Raquel Urtasun. Learning deep structured models. In *ICML*, 2015. 2, 3
- [12] Timothy F Cootes and Christopher J Taylor. Active shape models—‘smart snakes’. In *BMVC*, 1992. 1
- [13] Timothy F Cootes, Christopher J Taylor, David H Cooper, and Jim Graham. Active shape models-their training and application. *CVIU*, 1995.
- [14] Timothy F. Cootes, Gareth J. Edwards, and Christopher J. Taylor. Active appearance models. *IEEE TPAMI*, 2001.
- [15] David Cristinacce and Tim Cootes. Automatic feature localisation with constrained local models. *Pattern Recognition*, 2008. 1
- [16] Radek Danecek, Michael J. Black, and Timo Bolkart. EMOCA: Emotion driven monocular face capture and animation. In *Proc. CVPR*, 2022. 1
- [17] Ziqiang Dang, Jianfang Li, and Lin Liu. Cascaded dual vision transformer for accurate facial landmark detection. In *Proc. WACV*, 2025. 1
- [18] Arnaud Dapogny, Kevin Bailly, and Matthieu Cord. De-CaFA: Deep convolutional cascade for face alignment in the wild. In *Proc. ICCV*, 2019. 1, 2, 5, 6, 14
- [19] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. ArcFace: Additive angular margin loss for deep face recognition. In *Proc. CVPR*, 2019. 1
- [20] Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. In *CVPRW*, 2019. 1
- [21] Justin Domke. Structured learning via logistic regression. In *Proc. NeurIPS*, 2013. 2
- [22] Xuanyi Dong, Yan Yan, Wanli Ouyang, and Yi Yang. Style aggregated network for facial landmark detection. In *Proc. CVPR*, 2018. 1
- [23] Xuanyi Dong, Shou-I Yu, Xinshuo Weng, Shih-En Wei, Yi Yang, and Yaser Sheikh. Supervision-by-registration: An unsupervised approach to improve the precision of facial landmark detectors. In *Proc. CVPR*, 2018. 1
- [24] Xuanyi Dong, Yi Yang, Shih-En Wei, Xinshuo Weng, Yaser Sheikh, and Shou-I Yu. Supervision by registration and triangulation for landmark detection. *IEEE TPAMI*, 2020. 1
- [25] Yao Feng, Haiwen Feng, Michael J Black, and Timo Bolkart. Learning an animatable detailed 3d face model from in-the-wild images. *ACM ToG*, 2021. 1
- [26] Zhen-Hua Feng, Guosheng Hu, Josef Kittler, William Christmas, and Xiao-Jun Wu. Cascaded collaborative regression for robust facial landmark detection trained using a mixture of synthetic and real images with dynamic weighting. *IEEE TIP*, 2015. 1
- [27] Zhen-Hua Feng, Josef Kittler, William Christmas, Patrik Huber, and Xiao-Jun Wu. Dynamic attention-controlled cascaded shape regression exploiting training data augmentation and fuzzy-set sample weighting. In *Proc. CVPR*, 2017. 6, 15
- [28] Zhen-Hua Feng, Josef Kittler, Muhammad Awais, Patrik Huber, and Xiao-Jun Wu. Wing loss for robust facial landmark localisation with convolutional neural networks. In *Proc. CVPR*, 2018. 1, 2, 5, 6, 14, 15
- [29] Zheng Gao and Ioannis Patras. Self-supervised facial representation learning with facial region awareness. In *Proc. CVPR*, 2024. 1
- [30] Ross Girshick. Fast r-cnn. In *Proc. ICCV*, 2015. 4, 8
- [31] Colin Graber and Alexander Schwing. Graph structured prediction energy networks. In *Proc. NeurIPS*, 2019. 2
- [32] Colin Graber, Ofer Meshi, and Alexander Schwing. Deep structured prediction with nonlinear output transformations. In *Proc. NeurIPS*, 2018. 2
- [33] Shuyang Gu, Jianmin Bao, Hao Yang, Dong Chen, Fang Wen, and Lu Yuan. Mask-guided portrait editing with conditional gans. In *Proc. CVPR*, 2019. 1
- [34] Tong He, Zhi Zhang, Hang Zhang, Zhongyue Zhang, Junyuan Xie, and Mu Li. Bag of tricks for image classification with convolutional neural networks. In *Proc. CVPR*, 2019. 2
- [35] Yangyu Huang, Hao Yang, Chong Li, Jongyoo Kim, and Fangyun Wei. Adnet: Leveraging error-bias towards normal direction in face alignment. In *Proc. ICCV*, 2021. 1, 2, 3, 4, 5, 6, 12, 14, 15
- [36] Dalong Jiang, Yuxiao Hu, Shuicheng Yan, Lei Zhang, Hongjiang Zhang, and Wen Gao. Efficient 3d reconstruction for face recognition. *Pattern Recognition*, 2005. 1
- [37] Haibo Jin, Shengcai Liao, and Ling Shao. Pixel-in-pixel net: Towards efficient facial landmark detection in the wild. *IJCV*, 2021. 1, 5

- [38] Fatih Kahraman, Muhittin Gokmen, Sune Darkner, and Rasmus Larsen. An active illumination and appearance (AIA) model for face alignment. In *Proc. CVPR*, 2007. 1
- [39] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proc. ICLR*, 2015. 4, 12
- [40] Amit Kumar and Rama Chellappa. Disentangling 3d pose in a dendritic CNN for unconstrained 2d face alignment. In *Proc. CVPR*, 2018. 1
- [41] Abhinav Kumar, Tim K. Marks, Wenxuan Mou, Ye Wang, Michael Jones, Anoop Cherian, Toshiaki Koike-Akino, Xiaoming Liu, and Chen Feng. Luvli face alignment: Estimating landmarks’ location, uncertainty, and visibility likelihood. In *Proc. CVPR*, 2020. 1, 4, 5
- [42] John Lafferty, Andrew McCallum, Fernando Pereira, et al. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proc. ICML*, 2001. 2
- [43] Xing Lan, Qinghao Hu, and Jian Cheng. Revisting quantization error in face alignment. In *Proc. ICCVW*, 2021. 1, 3, 5, 6, 14, 15
- [44] Xing Lan, Jiayi Lyu, Kun Dong, Hanyu Jiang, Qinghao Hu, and Jian Xue. Does pixel value represent facial landmark well in heatmap? *IEEE TCSVT*, 2024. 2
- [45] Hui Li, Zidong Guo, Seon-Min Rhee, Seungju Han, and Jae-Joon Han. Towards accurate facial landmark detection via cascaded transformers. In *Proc. CVPR*, 2022. 2, 5, 6
- [46] Jinpeng Li, Haibo Jin, Shengcai Liao, Ling Shao, and Pheng-Ann Heng. Repformer: Refinement pyramid transformer for robust facial landmark detection. In *IJCAI*, 2022. 2, 5, 6
- [47] Weizhi Li, Gautam Dasarthy, and Visar Berisha. Regularization via structural label smoothing. In *Proc. AISTATS*, 2020. 2
- [48] Weijian Li, Yuhang Lu, Kang Zheng, Haofu Liao, Chihung Lin, Jiebo Luo, Chi-Tung Cheng, Jing Xiao, Le Lu, Chang-Fu Kuo, et al. Structured landmark detection via topology-adapting deep graph learning. In *Proc. ECCV*, 2020. 1, 5
- [49] Yujia Li and Rich Zemel. High order regularization for semi-supervised learning of structured output problems. In *Proc. ICML*. PMLR, 2014. 2
- [50] Yuanming Li, Gwantae Kim, Jeong-gi Kwak, Bon-hwa Ku, and Hanseok Ko. Towards multi-domain face landmark detection with synthetic data from diffusion model. In *Proc. ICASSP*, 2024. 1
- [51] Jiayi Liang, Haotian Liu, Hongteng Xu, and Dixin Luo. Generalizable face landmarking guided by conditional face warping. In *Proc. CVPR*, 2024. 1
- [52] Chunze Lin, Beier Zhu, Quan Wang, Renjie Liao, Chen Qian, Jiwen Lu, and Jie Zhou. Structure-coherent deep feature learning for robust face alignment. *IEEE TIP*, 2021. 1, 15
- [53] Guosheng Lin, Chunhua Shen, Ian Reid, and Anton van den Hengel. Deeply learning the messages in message passing inference. In *Proc. NeurIPS*, 2015. 2
- [54] Zhiwei Liu, Xiangyu Zhu, Guosheng Hu, Haiyun Guo, Ming Tang, Zhen Lei, Neil M Robertson, and Jinqiao Wang. Semantic alignment: Finding semantically consistent ground-truth for facial landmark detection. In *Proc. CVPR*, 2019. 1, 4
- [55] Iacopo Masi, Yue Wu, Tal Hassner, and Prem Natarajan. Deep face recognition: A survey. In *SIBGRAPI*, 2018. 1
- [56] Iain Matthews and Simon Baker. Active appearance models revisited. *IJCV*, 2004. 1
- [57] Paul Micaelli, Arash Vahdat, Hongxu Yin, Jan Kautz, and Pavlo Molchanov. Recurrence without recurrence: Stable video landmark detection with deep equilibrium models. In *Proc. CVPR*, 2023. 1
- [58] Stephen Milborrow and Fred Nicolls. Locating facial features with an extended active shape model. In *Proc. ECCV*, 2008. 1
- [59] Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. When does label smoothing help? In *Proc. NeurIPS*, 2019. 2
- [60] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *Proc. ECCV*, 2016. 2
- [61] Aiden Nibali, Zhen He, Stuart Morgan, and Luke Prendergast. Numerical coordinate regression with convolutional neural networks. *arXiv preprint arXiv:1801.07372*, 2018. 1, 3
- [62] Gabriel Pereyra, George Tucker, Jan Chorowski, Łukasz Kaiser, and Geoffrey Hinton. Regularizing neural networks by penalizing confident output distributions. In *Proc. ICLR*, 2017. 2
- [63] Shengju Qian, Keqiang Sun, Wayne Wu, Chen Qian, and Jiaya Jia. Aggregation via separation: Boosting facial landmark detector with semi-supervised style translation. In *Proc. ICCV*, 2019. 1
- [64] Shenhan Qian, Tobias Kirschstein, Liam Schoneveld, Davide Davoli, Simon Giebenhain, and Matthias Nießner. Gaussianavatars: Photorealistic head avatars with rigged 3d gaussians. In *Proc. CVPR*, 2024. 1
- [65] Joseph P Robinson, Yuncheng Li, Ning Zhang, Yun Fu, and Sergey Tulyakov. Laplace landmark localization. In *Proc. ICCV*, 2019. 2
- [66] Christos Sagonas, Georgios Tzimiropoulos, Stefanos Zafeiriou, and Maja Pantic. 300 faces in-the-wild challenge: The first facial landmark localization challenge. In *Proc. ICCVW*, 2013. 1, 6, 12
- [67] Jason Saragih and Roland Goecke. A nonlinear discriminative approach to aam fitting. In *ICCV*, 2007. 1
- [68] Patrick Sauer, Timothy F Cootes, Christopher J Taylor, et al. Accurate regression procedures for active appearance models. In *BMVC*, 2011. 1
- [69] Alexander Schwing and Raquel Urtasun. ICCV 2015 tutorial on learning deep structured models. https://www.cs.toronto.edu/~aschwing/ICCV_2015_DeepStructured_II_AlexSchwing.pdf, 2015. 2, 3, 12
- [70] Alexander G Schwing and Raquel Urtasun. Fully connected deep structured networks. *arXiv preprint arXiv:1503.02351*, 2015. 2
- [71] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Deep convolutional network cascade for facial point detection. In *Proc. CVPR*, 2013. 1

- [72] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proc. CVPR*, 2016. 2
- [73] Ying Tai, Yicong Liang, Xiaoming Liu, Lei Duan, Jilin Li, Chengjie Wang, Feiyue Huang, and Yu Chen. Towards highly accurate and stable face alignment for high-resolution videos. In *Proc. AAAI*, 2019. 2
- [74] Ben Taskar, Carlos Guestrin, and Daphne Koller. Max-margin markov networks. In *Proc. NeurIPS*, 2003. 2, 3, 12
- [75] Jonathan J Tompson, Arjun Jain, Yann LeCun, and Christoph Bregler. Joint training of a convolutional network and a graphical model for human pose estimation. In *Proc. NeurIPS*, 2014. 2
- [76] George Trigeorgis, Patrick Snape, Mihalisis A Nicolaou, Epameinondas Antonakos, and Stefanos Zafeiriou. Mnemonic descent method: A recurrent process applied for end-to-end face alignment. In *Proc. CVPR*, 2016. 2, 15
- [77] Ioannis Tsochantaridis, Thorsten Joachims, Thomas Hofmann, Yasemin Altun, and Yoram Singer. Large margin methods for structured and interdependent output variables. *JMLR*, 2005. 2, 3, 12
- [78] Roberto Valle, Jose M Buenaposada, Antonio Valdes, and Luis Baumela. A deeply-initialized coarse-to-fine ensemble of regression trees for face alignment. In *Proc. ECCV*, 2018. 6, 15
- [79] Jun Wan, Zhihui Lai, Jun Liu, Jie Zhou, and Can Gao. Robust face alignment by multi-order high-precision hourglass network. *IEEE TIP*, 2020. 6
- [80] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Minghui Tan, Xinggang Wang, Wenyu Liu, and Bin Xiao. Deep high-resolution representation learning for visual recognition. *IEEE TPAMI*, 2019. 1, 2, 5, 6
- [81] Xinyao Wang, Liefeng Bo, and Li Fuxin. Adaptive wing loss for robust face alignment via heatmap regression. In *Proc. ICCV*, 2019. 1, 4, 5, 6, 12, 14, 15
- [82] Erroll Wood, Tadas Baltrusaitis, Charlie Hewitt, Matthew Johnson, Jingjing Shen, Nikola Milosavljevic, Daniel Wilde, Stephan Garbin, Toby Sharp, Ivan Stojiljkovic, et al. 3d face reconstruction with dense landmarks. *arXiv*, 2022. 1
- [83] Wenyan Wu and Shuo Yang. Leveraging intra and inter-dataset variations for robust face alignment. In *Proc. CVPRW*, 2017. 15
- [84] Wayne Wu, Chen Qian, Shuo Yang, Quan Wang, Yici Cai, and Qiang Zhou. Look at boundary: A boundary-aware face alignment algorithm. In *Proc. CVPR*, 2018. 1, 2, 4, 6, 12, 14, 15
- [85] Yue Wu and Qiang Ji. Robust facial landmark detection under significant head poses and occlusion. In *Proc. ICCV*, 2015. 6, 15
- [86] Jiahao Xia, Weiwei Qu, Wenjian Huang, Jianguo Zhang, Xi Wang, and Min Xu. Sparse local patch transformer for robust face alignment and landmarks inherent relation learning. In *Proc. CVPR*, 2022. 1, 5, 6, 15
- [87] Shengtao Xiao, Jiashi Feng, Junliang Xing, Hanjiang Lai, Shuicheng Yan, and Ashraf Kassim. Robust facial landmark detection via recurrent attentive-refinement networks. In *Proc. ECCV*, 2016. 15
- [88] Xuehan Xiong and Fernando De la Torre. Supervised descent method and its applications to face alignment. In *Proc. CVPR*, 2013. 1, 15
- [89] Jing Yang, Qingshan Liu, and Kaihua Zhang. Stacked hourglass network for robust facial landmark localisation. In *Proc. CVPR*, 2017. 2
- [90] Raymond A. Yeh, Jinjun Xiong, Wen-Mei Hwu, Minh Do, and Alexander Schwing. Interpretable and globally optimal prediction for textual grounding using image concepts. In *Proc. NeurIPS*, 2017. 2
- [91] Baosheng Yu and Dacheng Tao. Heatmap regression via randomized rounding. *IEEE TPAMI*, 2021. 2
- [92] Hao Zhang, Lumin Xu, Shenqi Lai, Wenqi Shao, Nanning Zheng, Ping Luo, Yu Qiao, and Kaipeng Zhang. Open-vocabulary animal keypoint detection with semantic-feature matching. *IJCV*, 2024. 1
- [93] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang. Facial landmark detection by deep multi-task learning. In *Proc. ECCV*, 2014. 15
- [94] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang. Learning deep representation for face alignment with auxiliary attributes. *IEEE TPAMI*, 2015. 2
- [95] Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong Du, Chang Huang, and Philip HS Torr. Conditional random fields as recurrent neural networks. In *Proc. ICCV*, 2015. 2
- [96] Yinglin Zheng, Hao Yang, Ting Zhang, Jianmin Bao, Dongdong Chen, Yangyu Huang, Lu Yuan, Dong Chen, Ming Zeng, and Fang Wen. General facial representation learning in a visual-linguistic manner. In *Proc. CVPR*, 2021. 1
- [97] Erjin Zhou, Haoqiang Fan, Zhimin Cao, Yuning Jiang, and Qi Yin. Extensive facial landmark localization with coarse-to-fine convolutional network cascade. In *Proc. ICCVW*, 2013. 2
- [98] Zhenglin Zhou, Huaxia Li, Hong Liu, Nanyang Wang, Gang Yu, and Rongrong Ji. Star loss: Reducing semantic ambiguity in facial landmark detection. In *Proc. CVPR*, 2023. 1, 2, 3, 4, 5, 6, 7, 8, 12, 14, 15
- [99] Shizhan Zhu, Cheng Li, Chen Change Loy, and Xiaoou Tang. Face alignment by coarse-to-fine shape searching. In *Proc. CVPR*, 2015. 15
- [100] Xu Zou, Sheng Zhong, Luxin Yan, Xiangyun Zhao, Jiahuan Zhou, and Ying Wu. Learning robust facial landmark detection via hierarchical structured ensemble. In *Proc. CVPR*, 2019. 15

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.

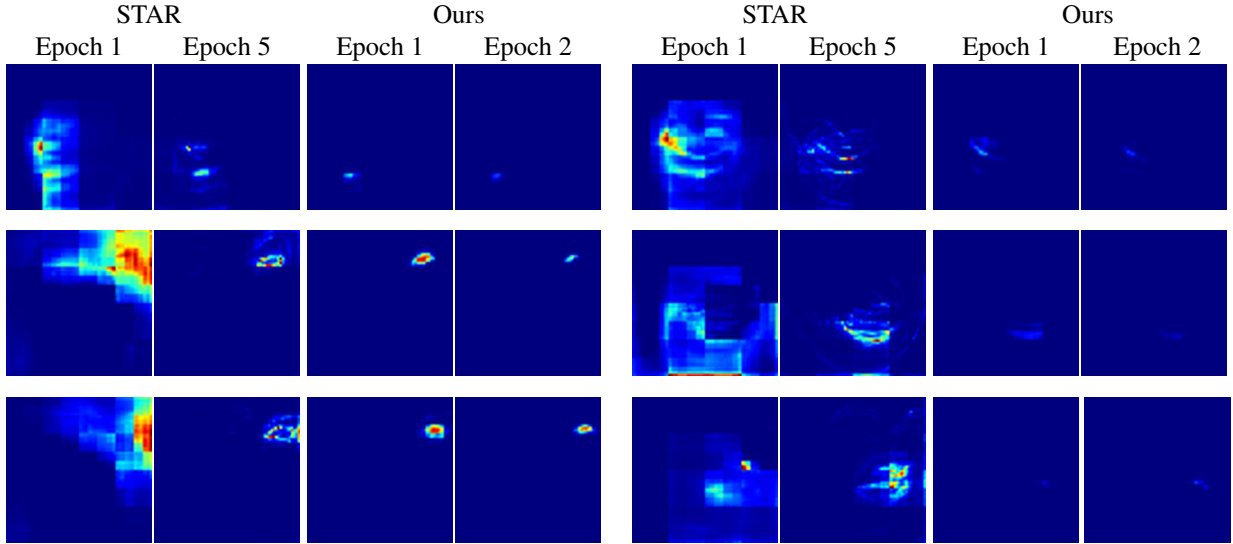


Figure A1. More comparisons of convergence efficiency between STAR’s Soft-argmax and ours following the settings in Fig. 4.

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 a 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 our assumption that the variation 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]. AAM 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 λ 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]

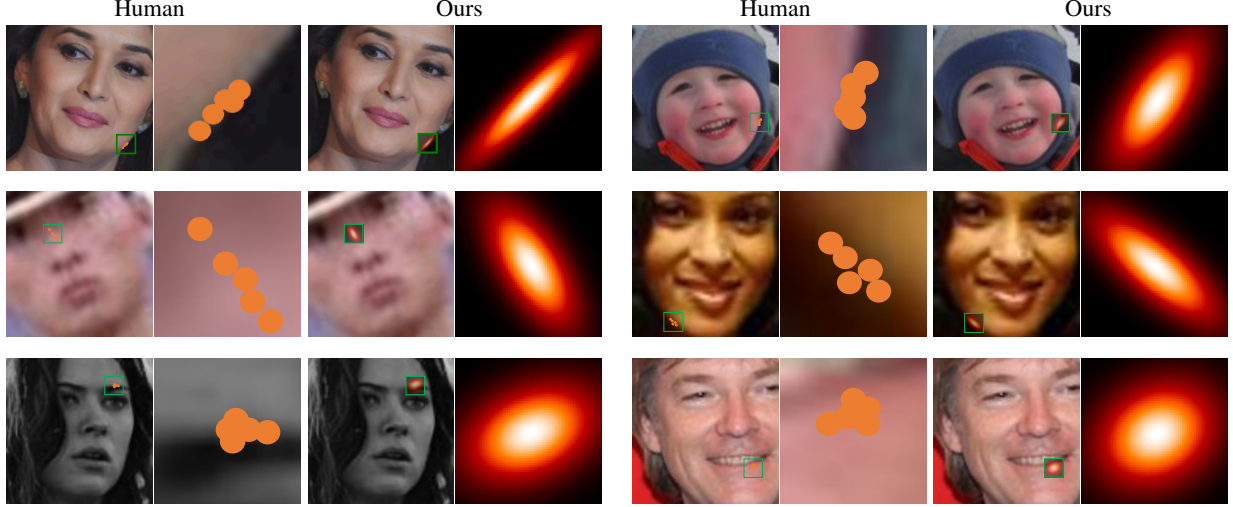


Figure A2. **Comparison between real-life annotation ambiguity and our label smoothing.** The “Human” columns are annotated by 5 users. The “Ours” columns are generated from our label smoothing.

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 γ 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: α (Eq. 12), ϵ (Eq. 11), and λ (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 ϵ and λ . Overall, we choose $\alpha = 1$, $\epsilon = 1$ and $\lambda = 5000$.

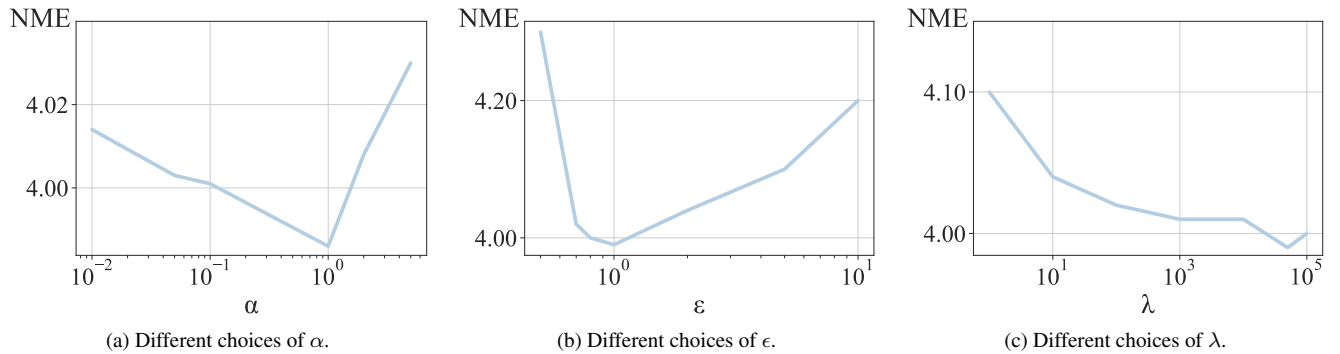


Figure A3. **The NME(↓) of different choices of hyper-parameters on WFLW.**

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	4.26	3.94	3.93	4.77	4.56
\checkmark	6.58	4.26	3.90	3.89	4.74	<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	1.27	2.43	<u>1.45</u>	5.16	3.10
STAR [98]	H	<u>10.77</u>	2.24	<u>1.58</u>	0.98	<u>4.76</u>	<u>2.98</u>
Ours	H	9.23	<u>1.92</u>	0.86	1.46	3.95	2.33

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	0.601	0.613	0.618	0.539	0.561
STAR [98]	H	<u>0.362</u>	0.584	0.609	0.622	<u>0.538</u>	0.551
Ours	H	0.370	<u>0.587</u>	<u>0.612</u>	<u>0.619</u>	0.539	<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	0.00
SLPT [86]	4.79	1.18	3.32	0.00
Awing [81]	4.94	0.99	-	-
ADNet [35]	4.68	<u>0.59</u>	-	-
HIH [43]	4.63	0.39	-	-
STAR [98]	<u>4.62</u>	0.79	<u>3.21</u>	0.00
Ours	4.58	0.79	3.15	0.00

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	3.27	7.18
Awing [81]	4.31	3.77	6.52
ADNet [35]	4.08	3.51	6.47
STAR [98]	4.03	3.55	6.22
Ours	4.03	3.51	6.09

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.

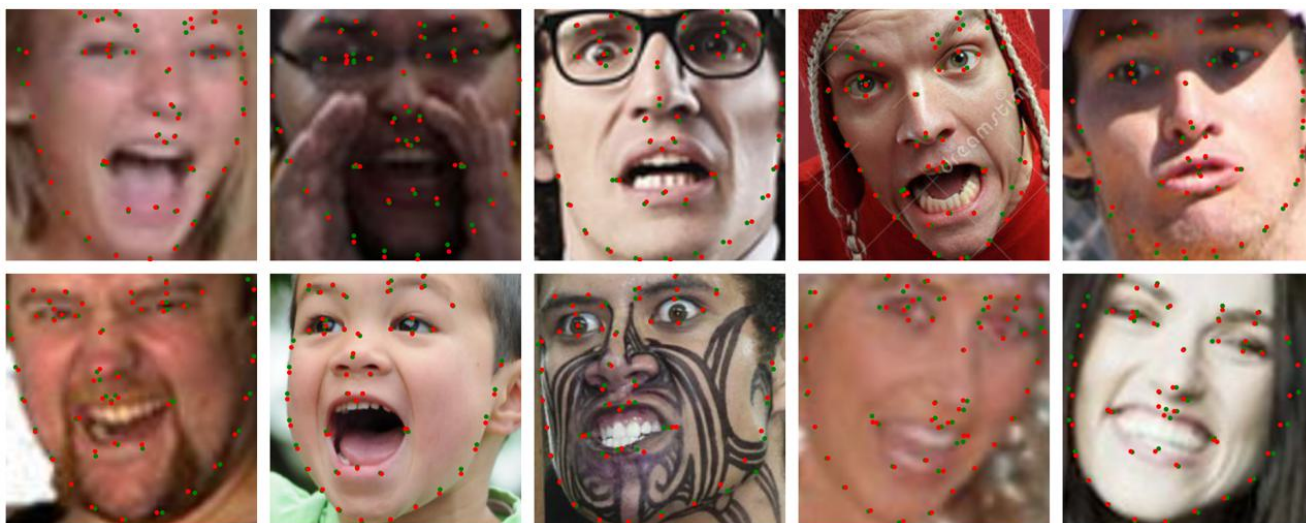


Figure A5. **Qualitative results on the “expression” subset of WFLW.** The ground truth landmarks are marked in **green** while our predictions are marked in **red**.



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.

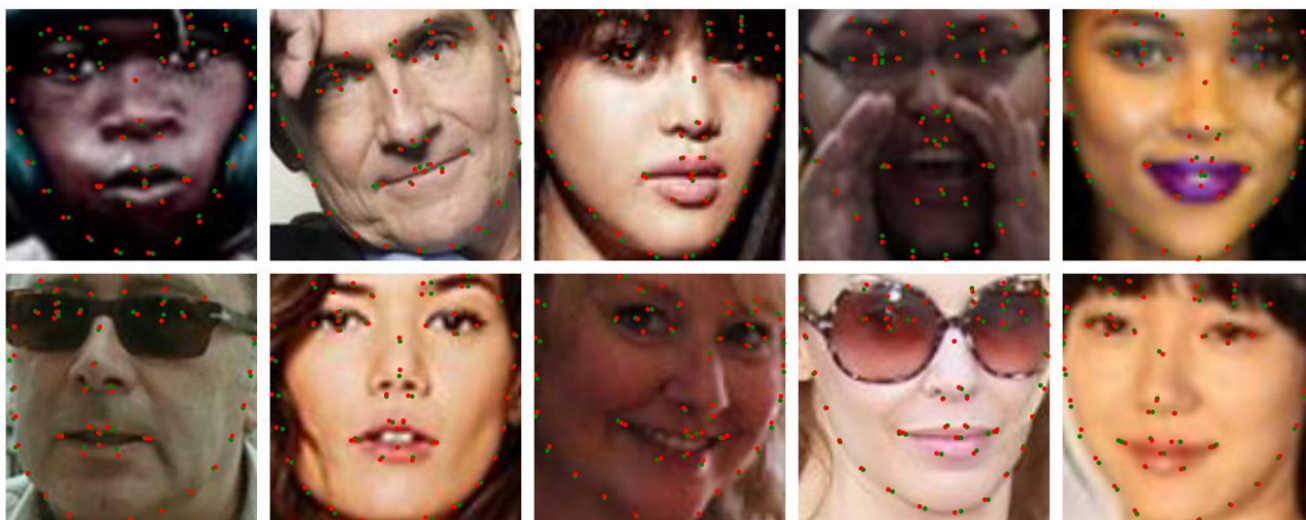


Figure A8. **Qualitative results on the “occlusion” subset of WFLW.** The ground truth landmarks are marked in green while our predictions are marked in red.

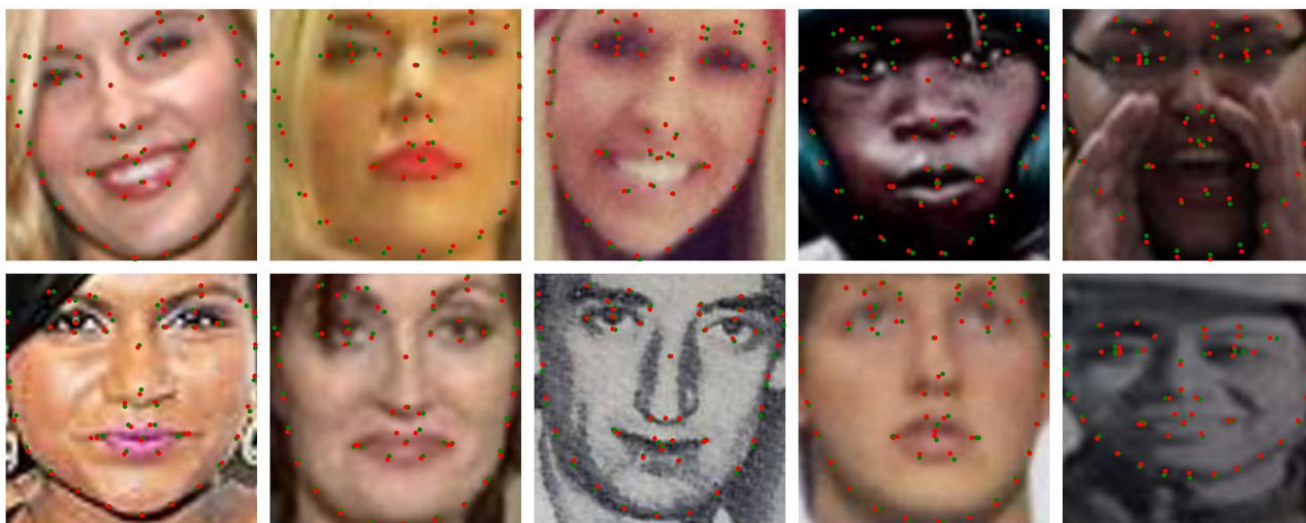


Figure A9. **Qualitative results on the “blur” subset of WFLW.** The ground truth landmarks are marked in **green** while our predictions are marked in **red**.



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**.