

3D-aware Facial Landmark Detection via Multiview Consistent Training on Synthetic Data

Anonymous CVPR submission

Paper ID 4794

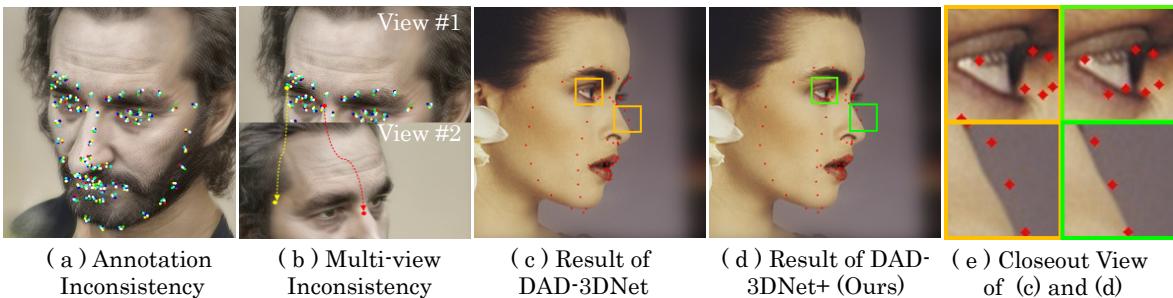


Figure 1. We plot the landmark annotations labeled by different annotators with different colors in (a). We project two annotated 3D points from View #1 to View #2, and we can find that the points are inaccurate if we look at the projected points in view #2. We show the landmark detection improvement in (c) (d), and (e). After refined by the proposed 3D-aware learning, the detected facial landmark is better aligned with the identity.

Abstract

Accurate facial landmark detection on wild images plays essential role for human-computer interaction, entertainment, and medical applications. Existing approaches are limited by the in-the-wild data and fail to enforce 3D consistency when detecting 3D/2D facial landmarks. On the other hand, with recent advances in generative visual models and neural rendering, we have witnessed rapid progress towards high quality 3D image synthesis. In this work, by leveraging synthetic data, we propose a novel multi-view consistent learning strategy to improve 3D facial landmark detection accuracy on in-the-wild images. The proposed 3D-aware module can be plugged into any learning-based landmark detection algorithm. We demonstrate the superiority of the proposed plug-in module with extensive comparison against state-of-the-art methods on several real and synthetic datasets.

1. Introduction

Accurate and precise facial landmark detection plays a significant role in computer vision and computer graphics

applications, such as face morphing [53], facial reenactment [57], 3D face reconstruction [17, 18, 30], head pose estimation [38], face recognition [1, 10, 13, 19, 32, 41, 69], face generation [11, 21, 59, 67], to name a few. In these applications, facial landmark detection provides great sparse representation to ease the burden of network convergence in different training stages. For instance, as a facial prior, it provides good initialization for subsequent training [65–67, 74], good intermediate representation to bridge the gap between different modalities for content generation [11, 27, 50, 77], loss terms which regularize the facial expression [11, 51], or evaluation metrics to quantitatively evaluate the facial motion quality [52, 71, 76].

The aforementioned applications require that facial landmark detection to be accurate even with significantly varied facial appearance under different identities, facial expressions, and extreme head poses. Tremendous efforts have been devoted to boost the performance [15, 22–24, 29, 34, 40, 55, 62, 72, 73, 75, 80, 82], which often rely on manually annotated large-scale lab-controlled or in-the-wild image datasets [4, 34], in order to handle arbitrary facial expressions, head poses, illumination, facial occlusions, etc.

However, even with high cost of human labeling, consis-

108 *tent and accurate* manual annotation of landmarks remains
109 challenging [22,23,34]. It is very difficult, if not impossible,
110 to force a person to annotate the facial landmark keypoints
111 at the same pixel locations for faces of different poses, let
112 along different annotators under different labeling environ-
113 ments. Such inconsistent annotations in training images are
114 often the killing factor to learn an accurate landmark local-
115 ization model. For example, Fig. 1(a) shows the inconsis-
116 tency of human annotations, which limits the training con-
117 vergence of previous methods. In Fig. 1(b), we project the
118 landmarks detected by a state-of-the-art algorithm [34] from
119 view #1 to view #2. We can observe obvious 3D inconsis-
120 tencies, which lead to poor landmark detection accuracy on
121 facial images with extreme head pose.
122

123 To mitigate this annotation inconsistency issue, we pro-
124 pose to learn facial landmark detection by enforcing multi-
125 view consistency during training. Given the images of the
126 same facial identity captured with different head poses, in-
127 stead of detecting facial landmark at each separate facial im-
128 age, we propose a multi-view consistency supervision to lo-
129 cate facial landmark in a holistic 3D-aware manner. To en-
130 force multi-view consistency, we introduce self-projection
131 consistency loss and multiview landmark loss in training.
132 We also propose an annotation generation procedure to
133 exploit the merits of lab-controlled data (*e.g.*, multi-view
134 images, consistent annotations) and in-the-wild data (*e.g.*,
135 wide range of facial expressions, identities). Thanks to the
136 synthetic data, our method does not rely on human anno-
137 tation to obtain the accurate facial landmark locations. There-
138 fore, it alleviates the notorious impact of learning from in-
139 accurate and inconsistent annotations.
140

141 We formulate our solution as a plug-in 3D aware module,
142 which can be incorporated into any facial landmark detec-
143 tor and can boost a pre-trained model with higher accuracy
144 and multi-view consistency. Extensive experiments demon-
145 strate the effectiveness. The main contributions of our work
146 are as follows:

- 147 • To the best of our knowledge, for the first time, we
148 show that our synthetic dataset combines the merits of
149 lab captured face image data (*e.g.* multi-view) and the
150 in-the-wild face image datasets (*e.g.* appearance diver-
151 sity). Compared with the traditional graphics pipeline,
152 the face images generated by our scheme is more re-
153 alistic, leading to better domain adaptability in the
154 trained models. Compared with the generative models,
155 the data generated by our scheme is more controllable
156 and multi-view consistent.
- 157 • We propose a novel 3D-aware optimization module,
158 which can be plugged into any learning-based facial
159 landmark detection methods. The proposed module is
160 able to improve the detection accuracy on unseen im-
161 ages beyond the training data by simply refining the

162 baseline models with our synthesized data.
163

- 164 • We demonstrate the performance improvements of our
165 module built on top multiple baseline methods on sim-
166 ulated dataset, lab-captured datasets, and in-the-wild
167 datasets.
168
- 169 • We collect a simulated multi-view face dataset, DAD-
170 3DHeads-Syn, without domain gap, and will release it
171 to facilitate future research.
172

2. Related Work

173 In this section, we review face landmark datasets and de-
174 tection algorithms that are most related to our approach. We
175 also provide a brief review of data simulation tools related
176 to our work.
177

2.1. Face Landmark Detection Dataset

179 **Lab-controlled dataset.** Datasets under “controlled”
180 conditions [8, 20, 36, 39, 46, 47, 63, 64, 70] typically col-
181 lect video/images from indoor scenarios with certain re-
182 strictions, *e.g.* pre-defined expressions, head poses, etc. For
183 example, FaceScape dataset [64] contains 938 individuals
184 and each with 20 expressions using a 68-camera array un-
185 der controlled illumination and positions. Thus, it contains
186 aligned and consistent multi-view images and facial land-
187 mark annotations. However, the identities, poses, and ex-
188 pressions are limited. In addition, the environment condi-
189 tions are fully controlled. These result in limited general-
190 ization capability of models trained on this dataset. More-
191 over, the annotation workflow of such dataset is expensive
192 and hard to scale.
193

194 **In-the-wild dataset.** The boom of internet image sharing
195 enabled many “in-the-wild” facial landmark datasets [3, 7,
196 32, 48, 83] that are collected from the webs to facilitate fa-
197 cial landmark detection research. However, manually an-
198notating facial landmarks on in-the-wild images is a time-
199 consuming process and not scalable. Zhu et al. [81] released
200 300W-LP by extending the original 300W dataset with syn-
201 thetic images with extreme pose through image profiling of
202 frontal pose images. However, the novel view images are
203 generated by simply applying rotation matrix on the orig-
204 inal images, which leads to limited view range and poor
205 image quality. Meanwhile, 300W-LP lacks diversity in face
206 appearance and expression because of the intrinsic limita-
207 tions of 300W. Recently, Martyniuk *et al.* [34] introduce a
208 new dataset, DAD-3DHeads by proposing a novel anno-
209 tation scheme, in which a annotator can see the texture ren-
210 dered onto the 3D mesh with respect to their fitting to ver-
211 ify that the facial landmark annotation results are accurate.
212 The proposed scheme addresses the problems exhibited by
213 existing labeling tools, such as “guessing” the positions of
214 the correct landmarks for invisible parts of the head, thus
215

enabling accurate annotations. DAD-3DHeads dataset contains 44,898 in-the-wild images, covering extreme facial expressions, poses, and challenging illuminations. However, even with the help from the proposed annotation scheme, the DAD-3DHeads still has some drawbacks. First, annotations across different annotators are inconsistent (see Fig. 1) and those noisy annotations could affect the training of the detection network. Second, since the depth is estimated by FLAME [33], annotation accuracy is limited by the FLAME model. Third, it lacks of multi-view images such that the 3D consistency of landmark detection methods cannot be investigated.

2.2. Data Simulation

Simulation [26, 28, 35, 42, 44, 45, 49, 58, 60, 61, 68] is a useful tool in situations where training data for learning-based methods is expensive to annotate or even hard to acquire. For example, Zeng *et al.* [68] and Richardson *et al.* [42] used 3D Morphable Model (3DMM) to render training data with different lighting condition, identity, expression, and texture basis elements for reconstructing detailed facial geometry. However, the simulated images by these approaches lack of realism have severe domain gap compared with real world captures, limiting the usage. Bak *et al.* [2] adapt synthetic data using a CycleGAN [79] with a regularization term for preserving identities. Ayush *et al.* [56] use the images and latent code generated by StyleGAN [79] to train a controllable portrait image generation model. However, it is hard to control the attribute consistencies of images simulated by generative models, which limits the usage of the generated datasets.

2.3. Face Landmark Detection Algorithms

Traditional facial landmark detection methods leverage either holistic facial appearance information [12], or the global facial shape patterns [31, 83]. They yield reasonable results for images captured in lab-controlled environments with frontal faces and good lighting, however the performance on most of in-the-wild images is inferior.

Recently, deep learning-based algorithms have made promising progress on 2D facial landmark localization [15, 22–24, 29, 34, 40, 55, 62, 72, 73, 75, 80, 82] in terms of robustness, generalizability, and accuracy. FAN [6] constructs, for the first time, a very strong baseline by combining a state-of-the-art architecture for landmark localization and train it on a very large yet synthetically expanded 2D facial landmark dataset. To address self-occlusion and large appearance variation, Zhu *et al.* [80] propose a Cascaded Convolutional Neural Network and Optimized Weighted Parameter Distance Cost Loss function to formulate the priority of 3DMM parameters during training instead of predicting facial landmark keypoints. To further address the problems of shape reconstruction and pose estimation simul-

Dataset Type	Lab-Controlled	In-the-wild	Ours
Examples			
In-the-wild	✗	✓	✓
Large Scale	✗	✓	✓
Balanced	✓	✗	✓
Multiview Consistent	✓	✗	✓
Annotation Consistent	✓	✗	✓
Scalable	✗	✗	✓

Figure 2. The feature comparison of different type of datasets.

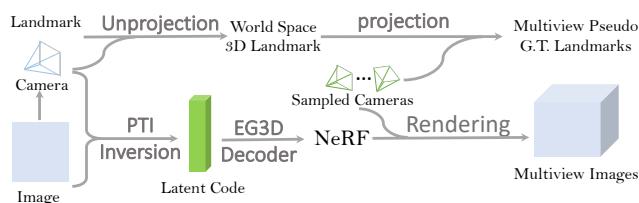


Figure 3. The proposed data simulation pipeline.

taneously, Martyniuk *et al.* propose an end-to-end trained DAD-3DNet [34] to regress 3DMM parameters and recovering the 3D head geometry with differential FLAME decoder. However, due to the intrinsic limitation of the manually annotated in-the-wild dataset, the detection results are affected by the annotation noise and the 3D inconsistency of the single view images. In this paper, we mainly focus on improving the performance of deep-learning based methods.

3. Balanced and Realistic Multi-view Faces Synthetic Dataset

Existing datasets are either lab-controlled captures [63, 64], or in-the-wild collected but lack of balanced multi-views [34, 81]. Moreover, we believe there are four desired properties that a good facial landmark dataset should fulfill: (1) containing full range of multi-view images; (2) can bridging the domain gap between the dataset and the real-world captured images; (3) containing diverse facial appearance including different poses, expressions, illuminations, and identities; (4) consistent and accurate annotations across the whole dataset (5) easy to obtain and scalable. Our dataset meets all the properties (Fig. 2).

Unlike previous graphics or generative model-based data synthesis approaches described in Sec. 1, we propose a novel facial dataset simulation scheme by leveraging Neural Radiance Field (NeRF) [37] to facilitate facial landmark detection tasks. Fig. 3 shows our dataset creation pipeline.

324 Specifically, we choose DAD-3DHeads [34] as our initial
 325 dataset since it contains images under a variety of extreme
 326 poses, facial expressions, challenging illuminations, and se-
 327 vere occlusions cases. As for the image generation, in-
 328 spired by GAN inversion [78], we first fit a latent code to
 329 each image in DAD-3DHeads datasets using EG3D [9] as
 330 decoder by following Pivotal Tuning Inversion (PTI) [43].
 331 Then we can use EG3D to decode it to NeRF. Next, we use
 332 volume rendering on the NeRF with 512 uniformly sam-
 333 pled camera views from a large view range, producing 512
 334 multi-view images. For landmark annotation, we start with
 335 the well-annotated groundtruth 3D landmarks of the origi-
 336 nal images from the DAD-3DHeads dataset. We first use
 337 Deep3DFace [14] to estimate the camera extrinsics with
 338 fixed camera intrinsics. Then we use the estimated camera
 339 view to unproject the annotated landmarks to 3D space. At
 340 last, we project the 3D landmarks to the 512 sampled cam-
 341 era views to obtain landmark annotation on the simulated
 342 views. The simulated dataset not only inherits the merits of
 343 DAD-3DHeads (*e.g.* diverse identities, expressions, poses,
 344 and illuminations), but also comes with a lot of new fea-
 345 tures (*e.g.*, balanced head pose, consistent annotation, and
 346 multi-view images). In total, there are 2,150,400 training
 347 pairs and 204,800 testing pairs in our extended dataset.
 348

4. 3D-Aware Multi-view Consistency Training

4.1. Overview

The state-of-the art landmark detectors [5,34] can output reasonable results on in-the-wild images. However, we may observe that the predicted landmark are floating on the face surface instead of fitting the face perfectly in a lot of cases. We can easily verify if the detected landmark fits the face by projecting the detected landmark to another view (see Fig. 1(b)). Innovated by this phenomenon, we propose a novel 3D-Aware training module \mathcal{R} to further improve the performance of baseline detection algorithm F .

Given a facial landmark detection network $F_\theta(\cdot)$ pre-trained on dataset \mathcal{D} , the proposed module \mathcal{R} further refine the network parameters θ by leveraging our simulated DAD-3DHead-Synth dataset $\hat{\mathcal{D}}$ and \mathcal{D} . Our module \mathcal{R} can be formulated as:

$$F_{\theta^*} \leftarrow \mathcal{R}(F_\theta, X, V_1, \dots, N), X \in \mathcal{D}, V_1, \dots, N \in \hat{\mathcal{D}}, \quad (1)$$

where X is the image batch sampled from \mathcal{D} and V_1, \dots, N are multi-view images sampled from $\hat{\mathcal{D}}$. Our goal is to refine the network parameters θ through exploring 3D information among multi-view images and applying a novel projection consistency during the fine-tuning process. Our module \mathcal{R} does not result in any new network parameters and can be plugged into any learning-based network. We show the training protocol in Alg. 1.

Algorithm 1 3D-Aware Plug-in Module.

Input: pretrained detector F with weights θ , M single-view images $I_1, \dots, M \in \mathcal{D}$ along with ground truth landmark L_1, \dots, M , paired N multi-view images $V_1, \dots, N \in \hat{\mathcal{D}}$ along with ground truth landmark L_1, \dots, N .
Output: detector F with updated weights θ^*
Initialization: set θ to pre-trained weights by it's initial training,
Unfreeze θ
for number of iterations **do**
 Output predicted landmarks \hat{L}_1, \dots, N for each view.
 Randomly sample P landmarks from them, $(1 < P \leq N)$.
 Cast the landmarks into world space and estimate the approximate 3D landmark \hat{L} using Eq. 2, 3, 4, 5
 Project \hat{L} onto the image planes of remaining Q views ($Q = N - P$) using Eq. 6, 7
 Calculate Self-Projection Consistency Loss $\mathcal{L}_{\text{Self-Cons}}$ using Eq. 8.
 Calculate Mesh Consistency Consistency Loss $\mathcal{L}_{\text{Mesh-Cons}}$ using Eq. 9
 Calculate Multiview Landmark Loss $\mathcal{L}_{\text{Multiview}}$ using Eq. 10
 Calculate Total Loss \mathcal{L} using Eq. 11
 $\theta^* \leftarrow Adam\{\mathcal{L}\}$

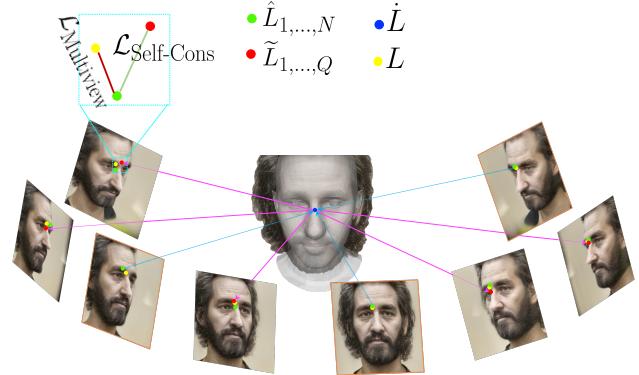


Figure 4. Multi-view Consistency Supervision. Predicted landmarks \hat{L}_1, \dots, N , estimated the 3D landmark \hat{L} , projection landmarks \tilde{L}_1, \dots, Q , and ground truth landmarks L are denoted as green, blue, red, and yellow points respectively. The process calculating 3D landmark \hat{L} and the projection procedure are shown as light blue and pink arrow lines respectively. $\mathcal{L}_{\text{Self-Cons}}$ and $\mathcal{L}_{\text{Multiview}}$ are represented as red line segment and light green segment respectively.

4.2. Multi-view Consistency Supervision

We propose a novel multi-view supervision to enforce the baseline network learn to be 3D consistent. To simplify notation, we ignore the batch dimension and fixed camera intrinsic matrix. For every training iteration, we randomly

sample N image and landmark pairs $\{V, L\}_{1,\dots,N}$ from \mathcal{D} and M image and landmark pairs $\{I, L\}_{1,\dots,M}$ from initial dataset \mathcal{D}^1 .

We pass $V_{1,\dots,N}$ to the baseline network F to obtain predicted landmarks $\hat{L}_{1,\dots,N}$ which are notated as green points in Fig. 4, and randomly select P predicted landmarks $\hat{L}_{1,\dots,P} \in \mathbb{R}^{P \times 68 \times 2}$ from $\hat{L}_{1,\dots,N}$ to calculate the “canonical” 3D landmark $\dot{L} \in \mathbb{R}^{68 \times 3}$ as shown by the blue point in Fig. 4 through Direct Linear Transformation (DLT) [16, 25] which is illustrated using blue arrowed lines in Fig. 4. Since all P views are simulated with pre-defined camera views using volume rendering (see Sec.3), the camera extrinsic matrix $\mathbb{M}_{1,\dots,P}$ are known. We calculate each keypoint of the “canonical” 3D landmark $\dot{L}^{(k)} \in \mathbb{R}^3, 1 \leq k \leq 68$, as follows:

$$\mu_p = \mathbb{M}_p[0, :] - \mathbb{M}_p[2, :] \cdot \hat{L}_p^k[0] \in \mathbb{R}^4, \quad (2)$$

$$v_p = \mathbb{M}_p[1, :] - \mathbb{M}_p[2, :] \cdot \hat{L}_p^k[1] \in \mathbb{R}^4, \quad (3)$$

$$\mathbf{A} = [\mu_1 | \mu_2 | \dots | \mu_p | v_1 | v_2 | \dots | v_p]^T \in \mathbb{R}^{2P \times 4}, \quad (4)$$

$$\dot{L}^{(k)} = \begin{pmatrix} \mathbf{A}[:, :3]^T & \mathbf{A}[:, :3] \end{pmatrix}^{-1} \mathbf{A}[:, :3]^T (-\mathbf{A}[:, 3]), \quad (5)$$

where, $p, 1 \leq p \leq P$ is the index of views. By Eq. 2 and Eq. 3, we first calculate the projection constraints for $\dot{L}^{(k)}$, i.e., $\mu_p[:, 3] \cdot \dot{L}^{(k)} + \mu_p[3] = 0$, where ‘ \cdot ’ indicates the dot product. Then we can stack all of the constraints into $\mathbf{A} \in \mathbb{R}^{2P \times 4}$ by Eq. 4. At last, we compute $\dot{L}^{(k)}$ with a least square approach (Eq. 5).

After obtaining the “canonical” 3D landmark \dot{L} , we project it onto the image planes of rest of $Q = N - P$ views to obtain the projected landmark $\tilde{L}_{1,\dots,Q}$, shown as red points in Fig. 4, by the following equations:

$$s = \mathbb{M}_q[:, :3] \dot{L}^{(k)} + \mathbb{M}_q[:, 3] \in \mathbb{R}^{3 \times 1}, \quad (6)$$

$$\tilde{L}_q^{(k)} = \begin{bmatrix} s[0]/s[2] \\ s[1]/s[2] \end{bmatrix} \in \mathbb{R}^{2 \times 1}, \quad (7)$$

where, in our case, $1 \leq q \leq Q$.

Self-Projection Consistency Loss. Since all M views are sampled from one NeRF with different camera views, the predicted landmarks $\hat{L}_{1,\dots,Q}$ and the projected landmarks $\tilde{L}_{1,\dots,Q}$ should be consistent, illustrated using light blue segment connecting red point and green point. The loss can be formulated as

$$\mathcal{L}_{\text{Self-Cons}} = \sum_{q=1}^Q \|\hat{L}_q - \tilde{L}_q\|_1. \quad (8)$$

¹ \mathcal{D} is DAD-3DHeads dataset when training DAD-3DNet and is AFLW2000-3D when training 3DDFA.

Mesh Consistency Consistency Loss² Besides the self-projection consistency, all the N views also share one mesh topology in canonical space. Therefore, we apply a mesh consistency loss in canonical space calculated by:

$$\mathcal{L}_{\text{Mesh-Cons}} = \sum_{n=1}^N \|\hat{M}_n - \dot{M}\|_2, \quad (9)$$

where \hat{M}_n is the predicted mesh of view n in canonical space, and \dot{M} is the ground truth mesh of the original reference image.

Multiview Landmark Loss. We also minimize the distance between predicted 2D facial landmarks and the corresponding multi-view ground truth landmarks we obtained in Sec.3, which are denoted as yellow points in Fig. 4. The loss can be formulated in the following equation:

$$\mathcal{L}_{\text{Multiview}} = \sum_{q=1}^N \|\hat{L}_q - L_q\|_1. \quad (10)$$

We also incorporate the original loss of the baseline method computed with image and landmark pairs $\{I, L\}_{1,\dots,M}$ from dataset \mathcal{D} to stabilize our 3D-aware training. The overall loss is:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{\text{Self-Cons}} + \lambda_2 \mathcal{L}_{\text{Mesh-Cons}} + \lambda_3 \mathcal{L}_{\text{Multiview}} + \mathcal{L}_{\text{original}}, \quad (11)$$

where $\lambda_{1,2,3}$ are hyper parameters that weights the balance between each components. We set $\lambda_{1,2,3}$ to 0.1 empirically.

Note that our training is a plug-in module and can be incorporated into any existing facial landmark detector easily. For different pretrained models, we just need to change $\mathcal{L}_{\text{original}}$, while other novel loss components calculated on our balanced synthetic dataset \mathcal{D} can be applied directly. we show this plug-in capability on top of different baseline methods (*e.g.*, DAD-3DNet [34] and 3DDFA [22]), and demonstrate that our 3D-aware training indeed improves their performance (see Sec. 5).

5. Experiments

5.1. Experimental Settings

Training Details. We implement our algorithm in Pytorch and adopt ADAM to optimize the baseline networks. We run our 3D-aware training for 100 epochs with a batch size of 4, and a learning rate of 1×10^{-4} on each baseline network.

Dataset. Besides DAD-3DHeads, we use another three datasets to conduct the evaluations.

- **DAD-3DHeads** [34] is the state-of-the-art in-the-wild 3D head dataset, which contains dense, accurate annotations, and diverse facial appearances. It consists of

²We can apply it depending on whether the baseline network outputs mesh. In our case, the 3DDFA [22] and DAD-3DNet [34] both do.



Figure 5. The visual results of Dlib [31], FAN [5], 3DDFA [22], our refined 3DDFA+, 3DDFA-V2, DAD-3DNet [34], and our refined DAD-3DNet on images randomly sampled from DAD-3DHeads [34] testing set. We show the enlarged error region (while box) in the middle row.

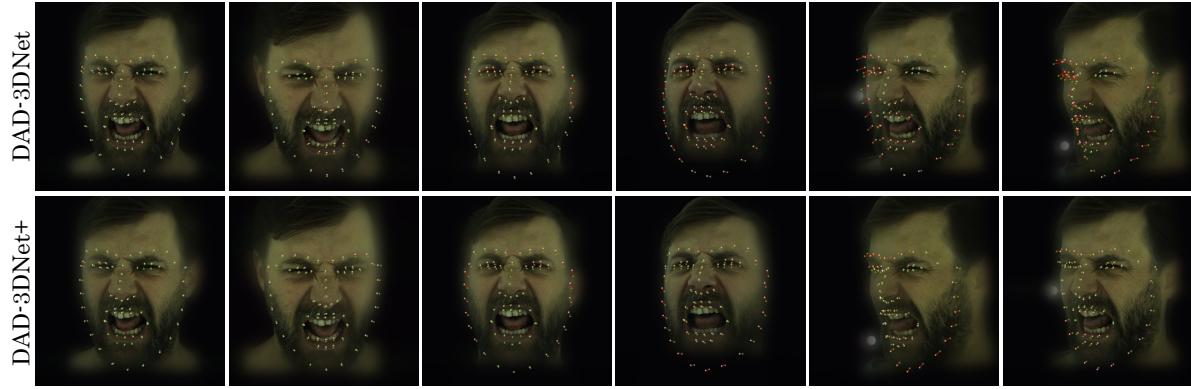


Figure 6. The error visualization of DAD-3DNet [34] and our DAD-3DNet+ on MultiFace [63] dataset. The white, green dots are the ground truth landmark and predicted landmark, respectively. We use the red line to show the error distance. From left to right, the head pose increases gradually.

Table 1. Facial landmark detection result (NME) on DAD-3DHeads [34], FaceScape [64], and MultiFace [63]. Lower values mean better results.

Method	DAD-3DHeads [34]	Multi-FaceScape [64]	Multi-Face [63]
FAN [6]	7.141	16.74	16.143
Dlib [31]	10.841	29.431	18.205
3DDFA-V2 [23]	2.926	6.853	5.942
3DDFA [22]	4.082	7.988	8.121
3DDFA+	3.784	7.425	7.305
DAD-3DNet [34]	2.599	6.681	5.786
DAD-3DNet+	2.503	6.050	5.480

44,898 images collected from various sources (37,840 in the training set, 4,312 in the validation set, and 2,746 in the test set).

- **FaceScape** [64] is a large-scale high-quality lab-controlled 3D face dataset, which contains 18,760 ex-

amples, captured from 938 subjects and each with 20 specific expressions.

- **MultiFace** [63] is a new multi-view, high-resolution human face dataset collected from 13 identities for neural face rendering.

Training and Testing Split. In all the experiments, we only refine the baseline models with the training set of our DAD-3DHeads-Syn and their original training dataset. We use the test sets of DAD-3DHeads-Syn and DAD-3DHeads [34], and use the full datasets of FaceScape [64] and MultiFace [62] for performance evaluation. All the comparison methods have not been trained on the split test sets.

Evaluation Metrics. We evaluate the facial landmark distance by calculating the normalized mean error (NME). We normalize the landmark error by dividing its image resolution instead of the eye distance [54], since all the test set images are aligned with offline tools. We calculate the head

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

Table 2. Head pose estimation results (head pose error) on DAD-3DHeads [34], FaceScape [64]. Lower values mean better results.

	DAD-3DHeads [34]				FaceScape [64]			
	Pitch	Roll	Yaw	Overall	Pitch	Roll	Yaw	Overall
FAN [5]	9.765	5.376	6.390	7.177	8.774	4.895	6.556	6.742
Dlib [31]	13.352	11.799	14.654	13.268	17.861	12.663	19.548	16.691
3DDFA-V2 [23]	7.901	4.989	6.088	6.326	13.741	9.718	11.353	11.604
3DDFA [22]	9.895	7.977	8.996	8.956	20.789	18.145	19.692	19.752
3DDFA+	9.195	6.792	8.692	8.226	20.996	16.426	19.054	18.826
DAD-3DNet [34]	8.274	4.666	9.206	7.382	15.851	9.676	18.346	14.624
DAD-3DNet+	7.700	4.274	7.528	6.500	14.466	7.247	13.876	11.863

pose error by the absolute distance of the Euler angle values.

5.2. Quantitative Evaluation

Landmark Detection Results. The quantitative results of landmark detection landmark detection results on DAD-3DHeads [34], FaceScape [64], and MultiFace [63] are shown in Tab.1. We can find that the DAD-3DNet+ refined by our 3D-aware multi-view consistency training achieves the best performance in all three datasets. Moreover, according to the results of 3DDFA [22], 3DDFA+, DAD-3DNet [34], and DAD-3DNet+, we find that after refinement, the new models (3DDFA+ and DAD-3DNet+) achieve much better results than the baseline models. For example, the detection error of DAD-3DNet [34] drops 0.631 and 0.306 on FaceScape and MultiFace datasets, respectively, which are 9% and 5% improvements. Similarly, we improve the 3DDFA [22] by 0.298 (7%), 0.563 (7%), and 0.816 (10%) on DAD-3DHeads, FaceScape and MultiFace datasets, respectively. And we attribute the improvement to our proposed 3D aware multi-view training. One interesting phenomenon is that all the methods perform better on DAD-3DHeads dataset than the other two lab-captured datasets. We attribute this to the extreme head pose and challenging facial expressions in the other two datasets. We plot the head pose distribution of DAD-3DHeads (see supplementary materials) and find that balance of head pose is not as good as other two lab-controlled datasets.

Head Pose Estimation Results. Table.2 shows the head pose estimation error on DAD-3DHeads [34] and FaceScape [64]. Our DAD-3DNet+ achieves best performance in most metrics. Similar as the landmark results, we can also conclude that our 3D-aware multi-view consistency training improves the head pose detection of the baseline methods (3DDFA and DAD-3DNet) based on the comparison results in forth row (3DDFA vs. 3DDFA+) and fifth row (DAD-3DNet vs. DAD-3DNet+). For example, after refinement, DAD-3DNet+ achieves 11.9%³ and 18.8%⁴ perfor-

³Head pose error drops from 7.382 to 6.500.

⁴Head pose error drops 14.624 to 11.863.

mance boosts in overall head pose error on DAD-3DHeads and FaceScape dataset, respetively.

5.3. Qualitative Evaluation

We fist show visual comparison on images randomly sampled from DAD-3DHeads testset [34] in Fig. 5. The landmark predicted by our DAD-3DNet+ model fits the individual’s face tighter than other predictions. Furthermore, By comparing third column (3DDFA [22]) and forth column (ours), we can find that the updated model our 3DDFA+ improve the landmark accuracy dramatically. Similar visual improvement can be found in sixth (DAD-3DNet) and seventh (DAD-3DNet+) columns as well. If we closely look at the landmark results (sixth and seventh column), we can find the refinement training can drag the landmark and rotate it in 3D space to better fit it with the individual’s face surface. We attribute this ability to our 3D-aware multi-view consistency training, which let the refined model gain the better sense in 3D space, therefor, improve the landmark detection results.

To further validate the improvement gained by the proposed 3D-aware multi-view consistency training, we show the visual results (Fig. 7) of 3DDFA [22], our refined 3DDFA+, DAD-3DNet [34], and our refined DAD-3DNet+ on images sampled from four different testsets. We can find that our proposed refinement improves the landmark detection results in eye, mouth, and face contour regions, which usually contain more appearance dynamics than other regions.

5.4. Performance Improvement Analysis

To systematically understand where is the improvement after refining the baseline methods (DAD-3DNet [34] and 3DDFA [22]) with our proposed 3D-aware multi-view consistency training, we further calculate and plot the landmark and head pose error improvements on DAD-3DHeads [34] (see Fig. 8). In stead of calculating the overall improved error score, we split all the testing images into different groups according to their head pose value and calculate the improved error score within each group. We can find that

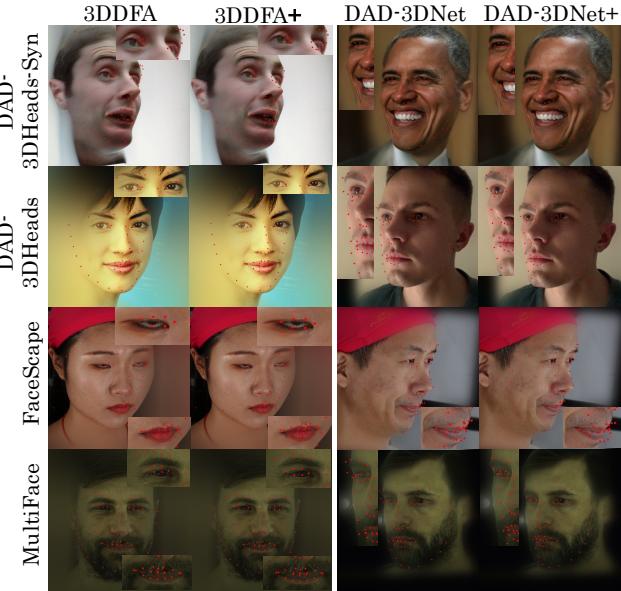


Figure 7. The visual comparisons between baseline methods and the refined methods on four testing sets. The left column and upper row list the dataset name and method name, respectively. '+' denotes the model that refined by our 3D-aware training.

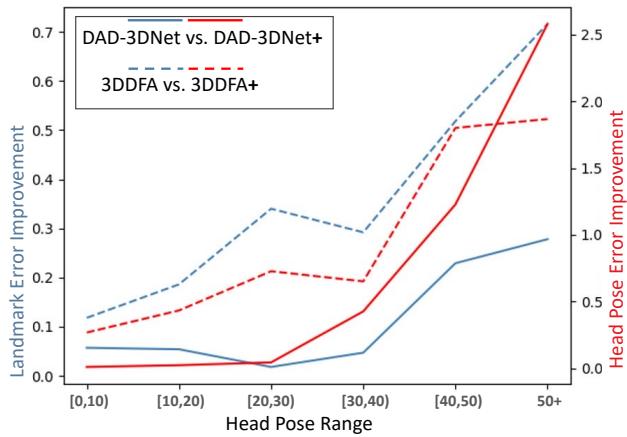


Figure 8. The landmark error and Head pose error improvement over DAD-3DNet [34] and 3DDFA [22] on images from different head pose ranges. The blue and red lines show the landmark error and head pose error improvement, respectively. The solid and dotted lines indicate DAD-3DNet [34] vs. DAD-3DNet+ (ours) and 3DDFA [22] vs. 3DDFA+ (ours).

the improvement by our training gets more obvious as the head pose gets more challenging. For example, the landmark error improvement (Fig. 8 blue part) using our method built on top of over 3DDFA [22] increases from 0.12 to 0.71. Similarly, the head pose estimation error (Fig. 6 ref part) improvement using our method built on top of DAD-3DNet [34] increases from 0.02 to 2.7. We also show the detection result visualization in Fig. 6. We can find that

Table 3. Ablation Study on FaceScape [64]

	Component	NME ↓	Pose ↓
1	full model (P=4)	6.050	11.863
2	w/o $\mathcal{L}_{\text{Mesh-Cons}}$	6.168	12.327
3	w/o $\mathcal{L}_{\text{Self-Cons}}$	6.541	13.623
4	full model (P=8)	6.048	11.923
5	full model (P=16)	6.098	11.902
6	full model (P=32)	6.139	11.912

from left to right, as the head pose increases, the error (ref line) of the DAD-3DNet+ (second row) is more stable than the error (first row) of the DAD-3DNet.

Base on this trend, we can conclude that our proposed 3D-aware multi-view consistency training improves more facial landmark detection and head pose estimation performance on images with larger head pose. This verifies our hypothesis that multi-view consistency training enable the network to learn 3D-aware information, therefor, benefiting the detection results on images with large head pose.

5.5. Ablation Study

We conduct ablation study on FaceScape [64] to verify the importance of main components of our novel design. As shown in Tab. 3, we calculate NME of landmark and MAE of pose estimation in these ablation experiments. Based on these numbers, we can see the performance degrades drastically when we remove $\mathcal{L}_{\text{Self-Cons}}$. Also, the role of $\mathcal{L}_{\text{Mesh-Cons}}$ is also verified crucial by the removal of it resulting sub-optimal result. As to the number of view used in the 3D triangulation procedure, it helps to generate good results. Moreover, using less views to do the estimation of 3D landmark in world space is more timing efficient and speedup the fine-tuning process.

6. Conclusion

We propose 3D-Aware Multi-View Consistency training, a new framework for improving deep-learning base landmark detection algorithms. We propose a novel Self-Projection Consistency Loss to better learn the 3D-aware information, a novel dataset simulation pipeline to combine the merits of lab-controlled captures and in-the-wild collected images. The model refined by our method outperforms previous approaches in terms of landmark detection accuracy and head pose estimation accuracy. Admittedly, our work has some limitations. For example, our proposed training relies on the performance of the baseline method. If the pretrianed baseline yield poor initial predictions, our DLT would fail to estimate reasonable canonical 3D landmark, affecting the performance of the proposed self-projection consistency loss.

864

References

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

- [1] Vitor Albiero, Xingyu Chen, Xi Yin, Guan Pang, and Tal Hassner. img2pose: Face alignment and detection via 6dof, face pose estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7617–7627, 2021.
- [2] Slawomir Bak, Peter Carr, and Jean-Francois Lalonde. Domain adaptation through synthesis for unsupervised person re-identification. In *Proceedings of the European conference on computer vision (ECCV)*, pages 189–205, 2018.
- [3] Peter N Belhumeur, David W Jacobs, David J Kriegman, and Neeraj Kumar. Localizing parts of faces using a consensus of exemplars. *IEEE transactions on pattern analysis and machine intelligence*, 35(12):2930–2940, 2013.
- [4] Adrian Bulat and Georgios Tzimiropoulos. Two-stage convolutional part heatmap regression for the 1st 3d face alignment in the wild (3dfaw) challenge. In *European Conference on Computer Vision*, pages 616–624. Springer, 2016.
- [5] Adrian Bulat and Georgios Tzimiropoulos. Binarized convolutional landmark localizers for human pose estimation and face alignment with limited resources. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3706–3714, 2017.
- [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 *Proceedings of the IEEE International Conference on Computer Vision*, pages 1021–1030, 2017.
- [7] Xavier P Burgos-Artizzu, Pietro Perona, and Piotr Dollár. Robust face landmark estimation under occlusion. In *Proceedings of the IEEE international conference on computer vision*, pages 1513–1520, 2013.
- [8] Chen Cao, Yanlin Weng, Stephen Lin, and Kun Zhou. 3d shape regression for real-time facial animation. *ACM Transactions on Graphics (TOG)*, 32(4):1–10, 2013.
- [9] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022.
- [10] Dong Chen, Shaoqing Ren, Yichen Wei, Xudong Cao, and Jian Sun. Joint cascade face detection and alignment. In *European conference on computer vision*, pages 109–122. Springer, 2014.
- [11] Lele Chen, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. Hierarchical cross-modal talking face generation with dynamic pixel-wise loss. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7832–7841, 2019.
- [12] Timothy F Cootes, Gareth J Edwards, and Christopher J Taylor. Active appearance models. In *European conference on computer vision*, pages 484–498. Springer, 1998.
- [13] Jiankang Deng, Jia Guo, Evangelos Ververas, Irene Kotsia, and Stefanos Zafeiriou. Retinaface: Single-shot multi-level face localisation in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5203–5212, 2020.

- [14] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019.
- [15] Xuanyi Dong, Yan Yan, Wanli Ouyang, and Yi Yang. Style aggregated network for facial landmark detection. In *proceedings of the IEEE conference on computer vision and pattern recognition*, June 2018.
- [16] Xuanyi Dong, Yi Yang, Shih-En Wei, Xinshuo Weng, Yaser Sheikh, and Shouo-I Yu. Supervision by registration and triangulation for landmark detection. *IEEE transactions on pattern analysis and machine intelligence*, 43(10):3681–3694, 2020.
- [17] Pengfei Dou, Shishir K Shah, and Ioannis A Kakadiaris. End-to-end 3d face reconstruction with deep neural networks. In *proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5908–5917, 2017.
- [18] Yao Feng, Fan Wu, Xiaohu Shao, Yanfeng Wang, and Xi Zhou. Joint 3d face reconstruction and dense alignment with position map regression network. In *Proceedings of the European conference on computer vision (ECCV)*, pages 534–551, 2018.
- [19] Golnaz Ghiasi and Charless C Fowlkes. Occlusion coherence: Detecting and localizing occluded faces. *arXiv preprint arXiv:1506.08347*, 2015.
- [20] Ralph Gross, Iain Matthews, Jeffrey Cohn, Takeo Kanade, and Simon Baker. Multi-pie. *Image and vision computing*, 28(5):807–813, 2010.
- [21] Kuangxiao Gu, Yuqian Zhou, and Thomas Huang. Flnet: Landmark driven fetching and learning network for faithful talking facial animation synthesis. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 10861–10868, 2020.
- [22] Jianzhu Guo, Xiangyu Zhu, and Zhen Lei. 3ddfa. <https://github.com/cleardusk/3DDFA>, 2018.
- [23] Jianzhu Guo, Xiangyu Zhu, Yang Yang, Fan Yang, Zhen Lei, and Stan Z Li. Towards fast, accurate and stable 3d dense face alignment. In *European Conference on Computer Vision*, pages 152–168. Springer, 2020.
- [24] Xiaojie Guo, Siyuan Li, Jinke Yu, Jiawan Zhang, Jiayi Ma, Lin Ma, Wei Liu, and Haibin Ling. Pfld: A practical facial landmark detector. *arXiv preprint arXiv:1902.10859*, 2019.
- [25] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003.
- [26] Stefan Hinterstoisser, Vincent Lepetit, Paul Wohlhart, and Kurt Konolige. On pre-trained image features and synthetic images for deep learning. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018.
- [27] Xinya Ji, Hang Zhou, Kaisiyuan Wang, Qianyi Wu, Wayne Wu, Feng Xu, and Xun Cao. Eamm: One-shot emotional talking face via audio-based emotion-aware motion model. In *ACM SIGGRAPH 2022 Conference Proceedings*, SIGGRAPH ’22, 2022.

- 972 [28] Justin Johnson, Bharath Hariharan, Laurens Van
973 Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross
974 Girshick. Clevr: A diagnostic dataset for compositional
975 language and elementary visual reasoning. In *Proceedings*
976 *of the IEEE conference on computer vision and pattern*
977 *recognition*, pages 2901–2910, 2017. 1026
- 978 [29] Amin Jourabloo and Xiaoming Liu. Large-pose face alignment
979 via cnn-based dense 3d model fitting. In *Proceedings of*
980 *the IEEE conference on computer vision and pattern recognition*,
981 pages 4188–4196, 2016. 1027
- 982 [30] Ira Kemelmacher-Shlizerman and Ronen Basri. 3d face re-
983 construction from a single image using a single reference
984 face shape. *IEEE transactions on pattern analysis and machine*
985 *intelligence*, 33(2):394–405, 2010. 1028
- 986 [31] Davis E. King. Dlib-ml: A machine learning toolkit. *Journal*
987 *of Machine Learning Research*, 10:1755–1758, 2009. 1029
- 988 [32] Martin Koestinger, Paul Wohlhart, Peter M Roth, and Horst
989 Bischof. Annotated facial landmarks in the wild: A large-
990 scale, real-world database for facial landmark localization.
991 In *2011 IEEE international conference on computer vision*
992 *workshops (ICCV workshops)*, pages 2144–2151. IEEE,
993 2011. 1030
- 994 [33] Tianye Li, Timo Bolkart, Michael. J. Black, Hao Li, and
995 Javier Romero. Learning a model of facial shape and ex-
996 pression from 4D scans. *ACM Transactions on Graphics*,
997 (*Proc. SIGGRAPH Asia*), 36(6), 2017. 1031
- 998 [34] Tetiana Martyniuk, Orest Kupyn, Yana Kurlyak, Igor
999 Krashenyi, Jiří Matas, and Viktoriia Sharmanska. Dad-
1000 3heads: A large-scale dense, accurate and diverse dataset
1001 for 3d head alignment from a single image. In *Proceedings*
1002 *of the IEEE conference on computer vision and pattern*
1003 *recognition*, pages 20942–20952, 2022. 1032
- 1004 [35] Nikolaus Mayer, Eddy Ilg, Philipp Fischer, Caner Hazir-
1005 bas, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox.
1006 What makes good synthetic training data for learning disparity
1007 and optical flow estimation? *International Journal of*
1008 *Computer Vision*, 126(9):942–960, 2018. 1033
- 1009 [36] Kieron Messer, Jiri Matas, Josef Kittler, Juergen Luettin,
1010 Gilbert Maitre, et al. Xm2vtsdb: The extended m2vts
1011 database. In *Second international conference on audio and*
1012 *video-based biometric person authentication*, volume 964,
1013 pages 965–966. Citeseer, 1999. 1034
- 1014 [37] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik,
1015 Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf:
1016 Representing scenes as neural radiance fields for view syn-
1017 thesis. *Communications of the ACM*, 65(1):99–106, 2021. 1035
- 1018 [38] Erik Murphy-Chutorian and Mohan Manubhai Trivedi. Head
1019 pose estimation in computer vision: A survey. *IEEE*
1020 *transactions on pattern analysis and machine intelligence*,
1021 31(4):607–626, 2008. 1036
- 1022 [39] P Jonathon Phillips, Patrick J Flynn, Todd Scruggs, Kevin W
1023 Bowyer, Jin Chang, Kevin Hoffman, Joe Marques, Jaesik
1024 Min, and William Worek. Overview of the face recogni-
1025 tion grand challenge. In *2005 IEEE computer society conference*
1026 *on computer vision and pattern recognition (CVPR'05)*, vol-
1027 ume 1, pages 947–954. IEEE, 2005. 1037
- 1028 [40] Shengju Qian, Keqiang Sun, Wayne Wu, Chen Qian, and Ji-
1029 aya Jia. Aggregation via separation: Boosting facial land-
1030 mark detector with semi-supervised style translation. In
1031 *Proceedings of the IEEE/CVF International Conference on*
1032 *Computer Vision*, pages 10153–10163, 2019. 1032
- 1033 [41] Rajeev Ranjan, Vishal M Patel, and Rama Chellappa. Hy-
1034 perface: A deep multi-task learning framework for face de-
1035 tection, landmark localization, pose estimation, and gender
1036 recognition. *IEEE transactions on pattern analysis and ma-*
1037 *chine intelligence*, 41(1):121–135, 2017. 1037
- 1038 [42] Elad Richardson, Matan Sela, Roy Or-El, and Ron Kimmel.
1039 Learning detailed face reconstruction from a single image.
1040 In *Proceedings of the IEEE conference on computer vision*
1041 and pattern recognition, pages 1259–1268, 2017. 1038
- 1042 [43] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel
1043 Cohen-Or. Pivotal tuning for latent-based editing of real im-
1044 ages. *ACM Transactions on Graphics (TOG)*, 42(1):1–13,
1045 2022. 1039
- 1046 [44] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Chris-
1047 tian Riess, Justus Thies, and Matthias Nießner. Faceforen-
1048 sics++: Learning to detect manipulated facial images. In
1049 *Proceedings of the IEEE/CVF international conference on*
1050 *computer vision*, pages 1–11, 2019. 1040
- 1051 [45] Nataniel Ruiz, Samuel Schulter, and Manmohan Chandraker.
1052 Learning to simulate. In *International Conference on Learn-
1053 ing Representations*, 2019. 1041
- 1054 [46] Christos Sagonas, Georgios Tzimiropoulos, Stefanos
1055 Zafeiriou, and Maja Pantic. 300 faces in-the-wild challenge:
1056 The first facial landmark localization challenge. In *Pro-
1057 ceedings of the IEEE international conference on computer*
1058 *vision workshops*, pages 397–403, 2013. 1042
- 1059 [47] Christos Sagonas, Georgios Tzimiropoulos, Stefanos
1060 Zafeiriou, and Maja Pantic. A semi-automatic methodology
1061 for facial landmark annotation. In *Proceedings of the IEEE*
1062 *conference on computer vision and pattern recognition*
1063 *workshops*, pages 896–903, 2013. 1043
- 1064 [48] Jie Shen, Stefanos Zafeiriou, Grigoris G Chrysos, Jean Kos-
1065 saifi, Georgios Tzimiropoulos, and Maja Pantic. The first
1066 facial landmark tracking in-the-wild challenge: Benchmark
1067 and results. In *Proceedings of the IEEE international con-
1068 ference on computer vision workshops*, pages 50–58, 2015. 1044
- 1069 [49] Ashish Shrivastava, Tomas Pfister, Oncel Tuzel, Joshua
1070 Susskind, Wenda Wang, and Russell Webb. Learning
1071 from simulated and unsupervised images through adversarial
1072 training. In *Proceedings of the IEEE conference on computer*
1073 *vision and pattern recognition*, pages 2107–2116, 2017. 1045
- 1074 [50] Linsen Song, Wayne Wu, Chaoyou Fu, Chen Change Loy,
1075 and Ran He. Audio-driven dubbing for user generated con-
1076 tents via style-aware semi-parametric synthesis. *IEEE Trans-
1077 actions on Circuits and Systems for Video Technology*, 2022. 1046
- 1078 [51] Linsen Song, Wayne Wu, Chen Qian, Ran He, and
1079 Chen Change Loy. Everybody's talkin': Let me talk as you
want. *IEEE Transactions on Information Forensics and Se-
curity*, 17:585–598, 2022. 1047
- 1080 [52] Yang Song, Jingwen Zhu, Dawei Li, Andy Wang, and
1081 Hairong Qi. Talking face generation by conditional recur-
1082 rent adversarial network. In *Proceedings of the Twenty-
1083 Eighth International Joint Conference on Artificial Intelli-
1084 gence, IJCAI-19*, pages 919–925. International Joint Confer-
1085 ences on Artificial Intelligence Organization, 7 2019. 1048

- 1080 [53] Luuk Spreeuwiers, Maikel Schils, and Raymond Veldhuis.
1081 Towards robust evaluation of face morphing detection. In
1082 *2018 26th European Signal Processing Conference (EUSIPCO)*, pages 1027–1031. IEEE, 2018.
1083
1084 [54] Keqiang Sun, Wayne Wu, Tinghao Liu, Shuo Yang, Quan
1085 Wang, Qiang Zhou, Zuochang Ye, and Chen Qian. Fab:
1086 A robust facial landmark detection framework for motion-
1087 blurred videos. In *Proceedings of the IEEE/CVF International
1088 Conference on Computer Vision*, pages 5462–5471,
1089 2019.
1090 [55] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Deep convolutional
1091 network cascade for facial point detection. In *Proceedings
1092 of the IEEE conference on computer vision and pattern
1093 recognition*, pages 3476–3483, 2013.
1094 [56] Ayush Tewari, Mohamed Elgharib, Gaurav Bharaj, Florian
1095 Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zoll-
1096 hofer, and Christian Theobalt. Stylerig: Rigging style-
1097 gan for 3d control over portrait images. In *Proceedings of
1098 the IEEE/CVF Conference on Computer Vision and Pattern
Recognition*, pages 6142–6151, 2020.
1099 [57] Justus Thies, Michael Zollhofer, Marc Stamminger, Chris-
1100 tian Theobalt, and Matthias Nießner. Face2face: Real-time
1101 face capture and reenactment of rgb videos. In *Proceed-
1102 ings of the IEEE conference on computer vision and pattern
1103 recognition*, pages 2387–2395, 2016.
1104 [58] Boris van Breugel, Trent Kyono, Jeroen Berrevoets, and Mi-
1105 haela van der Schaar. Decaf: Generating fair synthetic data
1106 using causally-aware generative networks. *Advances in Neu-
1107 ral Information Processing Systems*, 34:22221–22233, 2021.
1108 [59] Ting-Chun Wang, Ming-Yu Liu, Andrew Tao, Guilin Liu,
1109 Jan Kautz, and Bryan Catanzaro. Few-shot video-to-video
1110 synthesis. In *Advances in Neural Information Processing
1111 Systems (NeurIPS)*, 2019.
1112 [60] Erroll Wood, Tadas Baltrušaitis, Charlie Hewitt, Sebastian
1113 Dziadzio, Thomas J Cashman, and Jamie Shotton. Fake it
1114 till you make it: face analysis in the wild using synthetic
1115 data alone. In *Proceedings of the IEEE/CVF international
1116 conference on computer vision*, pages 3681–3691, 2021.
1117 [61] Erroll Wood, Tadas Baltrušaitis, Louis-Philippe Morency,
1118 Peter Robinson, and Andreas Bulling. Learning an
1119 appearance-based gaze estimator from one million synthe-
1120 sised images. In *Proceedings of the Ninth Biennial ACM
1121 Symposium on Eye Tracking Research & Applications*, pages
1122 131–138, 2016.
1123 [62] Yue Wu, Zuoguan Wang, and Qiang Ji. Facial feature track-
1124 ing under varying facial expressions and face poses based
1125 on restricted boltzmann machines. In *Proceedings of the
1126 IEEE Conference on Computer Vision and Pattern Rec-
ognition*, pages 3452–3459, 2013.
1127 [63] Cheng-hsin Wuu, Ningyuan Zheng, Scott Ardisson, Rohan
1128 Bali, Danielle Belko, Eric Brockmeyer, Lucas Evans, Timo-
1129 thy Godisart, Hyowon Ha, Alexander Hypes, Taylor Koska,
1130 Steven Krenn, Stephen Lombardi, Xiaomin Luo, Kevyn
1131 McPhail, Laura Millerschoen, Michal Perdoch, Mark Pitts,
1132 Alexander Richard, Jason Saragih, Junko Saragih, Takaaki
1133 Shiratori, Tomas Simon, Matt Stewart, Autumn Trimble,
Xinshuo Weng, David Whitewolf, Chenglei Wu, Shouo-I Yu,
and Yaser Sheikh. Multiface: A dataset for neural face ren-
1134 dering. In *arXiv*, 2022.
1135 [64] Haotian Yang, Hao Zhu, Yanru Wang, Mingkai Huang, Qiu
1136 Shen, Ruigang Yang, and Xun Cao. Facescape: a large-scale
1137 high quality 3d face dataset and detailed riggable 3d face pre-
1138 diction. In *Proceedings of the IEEE conference on computer
1139 vision and pattern recognition*, pages 601–610, 2020.
1140 [65] Ran Yi, Yong-Jin Liu, Yu-Kun Lai, and Paul L Rosin.
1141 Apdrawinggan: Generating artistic portrait drawings from
1142 face photos with hierarchical gans. In *Proceedings of the
1143 IEEE/CVF Conference on Computer Vision and Pattern
1144 Recognition*, pages 10743–10752, 2019.
1145 [66] Ran Yi, Zipeng Ye, Ruoyu Fan, Yezhi Shu, Yong-Jin Liu, Yu-
1146 Kun Lai, and Paul L Rosin. Animating portrait line drawings
1147 from a single face photo and a speech signal. In *ACM SIG-
GRAPH 2022 Conference Proceedings*, pages 1–8, 2022.
1148 [67] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and
Victor Lempitsky. Few-shot adversarial learning of realistic
1149 neural talking head models. In *Proceedings of the IEEE/CVF
1150 international conference on computer vision*, pages 9459–
1151 9468, 2019.
1152 [68] Xiaoxing Zeng, Xiaojiang Peng, and Yu Qiao. Df2net: A
1153 dense-fine-finer network for detailed 3d face reconstruction.
1154 In *Proceedings of the IEEE/CVF International Conference
1155 on Computer Vision*, pages 2315–2324, 2019.
1156 [69] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao.
1157 Joint face detection and alignment using multitask cascaded
1158 convolutional networks. *IEEE signal processing letters*,
1159 23(10):1499–1503, 2016.
1160 [70] Xing Zhang, Lijun Yin, Jeffrey F Cohn, Shaun Cana-
1161 van, Michael Reale, Andy Horowitz, and Peng Liu. A
1162 high-resolution spontaneous 3d dynamic facial expression
1163 database. In *2013 10th IEEE international conference and
1164 workshops on automatic face and gesture recognition (FG)*,
1165 pages 1–6. IEEE, 2013.
1166 [71] Zhimeng Zhang, Lincheng Li, Yu Ding, and Changjie
1167 Fan. Flow-guided one-shot talking face generation with
1168 a high-resolution audio-visual dataset. In *Proceedings of
1169 the IEEE/CVF Conference on Computer Vision and Pattern
1170 Recognition*, pages 3661–3670, 2021.
1171 [72] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou
1172 Tang. Facial landmark detection by deep multi-task learning.
1173 In *European conference on computer vision*, pages 94–108.
Springer, 2014.
1174 [73] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou
1175 Tang. Learning deep representation for face alignment with
auxiliary attributes. *IEEE transactions on pattern analysis
1176 and machine intelligence*, 38(5):918–930, 2015.
1177 [74] Aihua Zheng, Feixia Zhu, Hao Zhu, Mandi Luo, and Ran He.
1178 Talking face generation via learning semantic and temporal
1179 synchronous landmarks. In *2020 25th International Con-
ference on Pattern Recognition (ICPR)*, pages 3682–3689.
IEEE, 2021.
1180 [75] Erjin Zhou, Haoqiang Fan, Zhimin Cao, Yuning Jiang, and
Qi Yin. Extensive facial landmark localization with coarse-
1181 to-fine convolutional network cascade. In *Proceedings of
1182 the IEEE international conference on computer vision work-
1183 shops*, pages 386–391, 2013.
1184
1185
1186
1187

- 1188 [76] Hang Zhou, Yasheng Sun, Wayne Wu, Chen Change Loy, 1242
1189 Xiaogang Wang, and Ziwei Liu. Pose-controllable talking 1243
1190 face generation by implicitly modularized audio-visual rep- 1244
1191 resentation. In *Proceedings of the IEEE/CVF conference on* 1245
1192 *computer vision and pattern recognition*, pages 4176–4186, 1246
1193 2021. 1247
- 1194 [77] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevar- 1248
1195 ria, Evangelos Kalogerakis, and Dingzeyu Li. Makettalk: 1249
1196 speaker-aware talking-head animation. *ACM Transactions 1250
1197 on Graphics (TOG)*, 39(6):1–15, 2020. 1251
- 1198 [78] Jiapeng Zhu, Yujun Shen, Deli Zhao, and Bolei Zhou. In- 1252
1199 domain gan inversion for real image editing. In *European 1253
1200 conference on computer vision*, pages 592–608. Springer, 1254
2020. 1255
- 1201 [79] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A 1256
1202 Efros. Unpaired image-to-image translation using cycle- 1257
1203 consistent adversarial networks. In *Proceedings of the IEEE 1258
1204 international conference on computer vision*, pages 2223– 1259
1205 2232, 2017. 1260
- 1206 [80] Xiangyu Zhu, Zhen Lei, Xiaoming Liu, Hailin Shi, and 1261
1207 Stan Z Li. Face alignment across large poses: A 3d solu- 1262
1208 tion. In *Proceedings of the IEEE conference on computer 1263
1209 vision and pattern recognition*, pages 146–155, 2016. 1264
- 1210 [81] Xiangyu Zhu, Zhen Lei, Xiaoming Liu, Hailin Shi, and 1265
1211 Stan Z Li. Face alignment across large poses: A 3d solu- 1266
1212 tion. In *Proceedings of the IEEE conference on computer 1267
1213 vision and pattern recognition*, pages 146–155, 2016. 1268
- 1214 [82] Xiangyu Zhu, Xiaoming Liu, Zhen Lei, and Stan Z Li. Face 1269
1215 alignment in full pose range: A 3d total solution. *IEEE 1270
1216 transactions on pattern analysis and machine intelligence*, 1271
41(1):78–92, 2017. 1272
- 1217 [83] Xiangxin Zhu and Deva Ramanan. Face detection, pose es- 1273
1218 timation, and landmark localization in the wild. In *2012 1274
1219 IEEE conference on computer vision and pattern recogni- 1275
1220 tion*, pages 2879–2886. IEEE, 2012. 1276
- 1221 1277
- 1222 1278
- 1223 1279
- 1224 1280
- 1225 1281
- 1226 1282
- 1227 1283
- 1228 1284
- 1229 1285
- 1230 1286
- 1231 1287
- 1232 1288
- 1233 1289
- 1234 1290
- 1235 1291
- 1236 1292
- 1237 1293
- 1238 1294
- 1239 1295
- 1240
- 1241