

Survey on Face Tracking with Deep Learning

Vinay Balasubramanian¹ and Jilliam Diaz Barros²

¹ v.balasubr18@cs.uni-kl.de

² jilliam.maria.diaz.barros@dfki.de

Abstract. Face tracking is the underlying task for many applications such as face recognition, expression analysis, 3D face modeling, etc. In spite of recent progress with various state of the art methods, the task of detecting facial landmarks remain challenging in wild(unconstrained) conditions. In this paper, we review different face tracking architectures and their performance in challenging conditions. We focus on deep-learning based methods that exploit the temporal information across frames, i.e video-based methods. Recent developments include using an encoder-decoder network, recurrent network, deep reinforcement learning, and two-stream network. This paper aims to compare those approaches in terms of accuracy, the dataset(s) used for training, evaluation metrics, robustness to large head poses and occlusions, etc

Keywords: Face tracking, Facial landmarks, Deep Learning, Reinforcement Learning, Temporal information

1 Introduction

Face tracking is a computer vision task of tracking the face across all frames of a video. It may involve tracking specific landmarks around the face, or tracking a bounding box around the face across frames. Facial landmarks are mainly localized around facial components such as eyes, ears, nose, mouth and jawline. Face Tracking technology plays an important role in computer vision applications such as *Face recognition* [6], *Expression recognition* [3] and *Face modeling* [14]. This is a challenging problem as the videos may not be captured in constrained conditions and may have illumination inconsistencies, large head poses, blurriness, occlusions etc.

There are various approaches to this problem. Some of them are image-based methods, where the models are trained on still frames and the detection also happens independently at each frame. Other methods are video-based and use an incremental-learning technique to exploit the temporal connection between successive frames. Figure 1 shows a generic high-level architecture of a video-based landmark detection pipeline.

In this paper we make a comparison between the different architectures, datasets used for training and testing, evaluation metrics and robustness to challenging conditions.

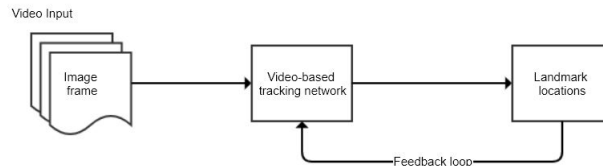


Fig. 1: Generic architecture of video based methods. Landmarks detected in the current frame are used as an initialization for the next frame

2 Datasets

In this section, we list the datasets commonly used for landmark-based face tracking. These datasets are publicly available for research purposes. Datasets can be categorized into constrained datasets and unconstrained datasets (in the wild). Table 1 shows various image-based and video datasets. Most methods use the 300-VW [18] and the TF [1] datasets to evaluate their models and compete in the 300-VW challenge.

Table 1: Datasets used for facial landmark detection and tracking

Dataset	Description	Video/Image	Contains	Wild?
AFLW [12]	Annotated Facial Landmarks in the Wild	Image	Around 25k annotated face images with 21 landmarks per image	Yes
COFW [5]	Caltech Occluded Faces in the Wild	Image	1007 occluded face images with 29 manually annotated landmarks	Yes
Helen [13]	Helen facial feature dataset	Image	2000 training and 330 test images with 194 landmarks and accurate annotations of primary facial components	Yes
IBUG [2]	IBUG dataset	Image	135 images with difficult poses and expressions	Yes
LFPW [4]	Labeled Face Parts in the Wild	Image	1432 images with 29 landmarks on each image	Yes
LFW [11]	Labeled Faces in the Wild	Image	13,233 images of 5749 people detected and centered by Viola Jones face detector	Yes
3D Menpo [20]	The 3D Menpo database	Image	14000 static images with 2D and 3D landmarks and around 280,000 annotated frames	Yes
BIWI [7]	Biwi kinect head pose database	Video	24 videos with over 15k frames of 20 people	Yes
FM [17]	Face Movies	Video	2150 images of 6 videos with 68 landmarks on each image	Yes
RWMB [19]	Real-World Motion Blur	Video	10000 face videos with 98 landmarks including occlusion, blur, illumination changes etc.	Yes
SynHead [8]	Synthetic dataset	Video	510,960 frames of 70 head motion tracks that include large face pose variations	Yes
TF [1]	Talking Face	Video	5000 frames of a person engaged in a conversation with 68 landmarks in each frame on each image	No
300VW [18],	300 videos in the wild	Video	114 videos with 218,595 frames with 68 landmarks per frame	Yes

3 Face Tracking Approaches

In this section, we describe some of the state-of-the-art approaches for video-based facial landmark tracking. Deep learning methods, in general, use CNN and RNN to detect landmarks.

3.1 Recurrent Encoder-Decoder Network for Video-based Face Alignment (2018) [16]

This method leverages temporal information to predict facial landmarks in each frame and uses recurrent learning at both spatial and temporal dimensions. At the temporal level, the features are separated into *temporal-variant* features such as pose and expression, and *temporal-invariant* features such as facial identity. Recurrent learning is only applied to the temporal-variant features. This feature disentangling has shown to achieve better generalization and more accurate results. Figure 2 shows the pipeline of recurrent encoder-decoder network.

The network consists of 4 modules -

- (1) **Encoder-Decoder**:- The encoder encodes features from a single video frame into an intermediate low dimensional representation by performing a sequence of convolutions, pooling and batch normalization. The decoder upsamples the low dimensional representation and transforms it into a response map that contains facial landmarks.

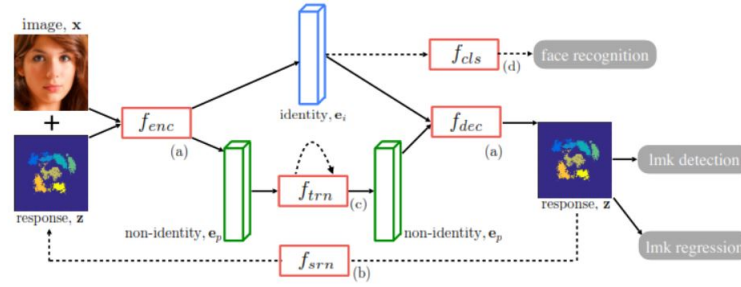


Fig. 2: Overview of REDNet pipeline Source: [16]

- (2) **Spatial recurrent learning**:- The purpose is to find the exact location of landmarks in a coarse-to-fine manner by iteratively providing the previous prediction as feedback along with the video frame. This is carried out in 2 steps - *Landmark Detection* and *Landmark Regression*. Landmark detection step locates 7 major facial components whereas landmark regression step refines predicted locations of all 68 landmark positions
- (3) **Temporal recurrent learning**:- This is proposed to model the temporal-variant factors such as pose and expression. The temporal variations in the temporal-invariant factors (non-identity code) are modeled using an LSTM unit consisting of 256 hidden neurons. Trained using T successive frames. Detection and regression tasks are performed frame by frame.
- (4) **Supervised identity disentangling**:- Complete identity and non-identity factor disentangling cannot be guaranteed. More supervised information is needed to achieve better separation of the features. This module applies identity constraint to the identity code to further separate identity code from the non-identity code. Face recognition is applied to the identity code to classify the people present in the frames. This is shown to yield better generalization and better test accuracy

3.2 Dynamic Facial Analysis using Recurrent Neural Networks (2017) [8]

This approach uses RNN for head pose estimation and facial landmark tracking. It proposes RNN as an alternative approach that performs better than previous video-based approaches for dynamic facial analysis which use Kalman filters or particle filters. The method is inspired by the fact that RNNs and Bayesian filters are operationally very similar although Bayesian filters need problem-specific hand-tuning. Given sufficient data, an RNN can be trained to do the same task and avoid problem-specific tracker engineering. The authors create a synthetic dataset **SynHead** to cater to the need for large training data.

The approach employs FC-RNN to exploit the generalization from a pre-trained CNN and consists of CNN layers followed by recurrent layers as dense layers. Figure 3 shows the proposed architecture for head pose estimation and tracking. The CNN and RNN are trained together end-to-end. The head pose is estimated in terms of pitch, yaw and roll angles. The network is a modified VGG16 with an extra fully connected layer with 1024 neurons and the output layer consists of 3 neurons for the pitch, yaw and roll angles. For facial landmark detection, the same network is used with the only difference that the output layer contains 136 neurons corresponding to the locations of the 68 landmarks. RNN makes the model robust to occlusions and large head poses.

3.3 Dual-Agent Deep Reinforcement Learning (2018) [9]

This approach exploits the fact that bounding box tracking and landmark detection are dependent. The accuracy of the detected facial landmarks depends on how good the bounding box is. Figure 4



Fig. 3: Proposed end-to-end CNN RNN network. Source: [8]

shows different strategies for deformable face tracking including the proposed DADRL (Dual-Agent Deep Learning) architecture. This framework is designed for simultaneous bounding box tracking and landmark detection in an interactive manner and uses reinforcement learning to learn to make adaptive decisions during face tracking. The architecture consists of a *Tracking agent*, an *Alignment agent* and *communication channels* between the agents. The two agents are trained simultaneously to learn two conditional distributions. Figure 5 shows the proposed architecture. The message channels are trained using deep Q-learning algorithm

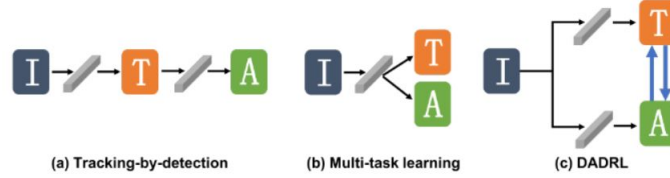


Fig. 4: Strategies for deformable face tracking. Source: [9]

If I_k is the k^{th} frame, B_k is the bounding box for the k^{th} frame and V_k is the vector of L landmarks, then the probabilistic duality is given by -

$$p(B_k|I_k) p(V_k|B_k, I_k) = p(V_k|I_k) p(B_k|V_k, I_k) \quad (1)$$

The learning objectives of bounding box tracking and landmark detection are treated as two conditional probabilities and the dependency between these two tasks is formulated as two marginal distributions. Since the ground-truth marginal distributions are not available, communication channels between the agents are used as alternatives to satisfy the probabilistic duality. For each frame, the terminal state of the previous frame is used for initializing the current state. The two agents decide a sequence of actions based on the observed state and exchanged messages, to adjust the bounding box and regress facial landmarks simultaneously. The messages sent from the tracking agent to the alignment agent are encoded by a deconvolution layer. It provides additional textural information to the alignment agent to improve its robustness. The messages from the alignment agent to the tracking agent are encoded by an LSTM unit. It provides 3D pose information to the tracking agent to improve bounding box tracking.

3.4 Two Stream Transformer Networks (2017) [15]

This approach proposes a two-stream deep learning method that decomposes the video input to spatial and temporal streams. The spatial stream aims to capture appearance information from still

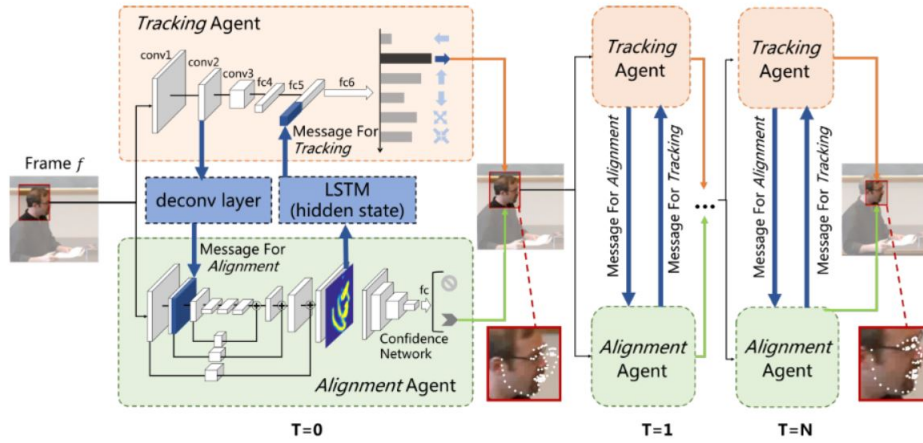


Fig. 5: DADRL architecture. Source: [9]

frames and it is trained to transform image pixels to landmark positions directly on still frames and then to refine the current facial shape based on the previous shape. On the other hand, the temporal stream aims to capture temporal consistency information across successive frames.

Figure 6 shows the proposed architecture. The temporal stream consists of an encoder-decoder module. The encoder is trained to encode the spatial information as active appearance codes that capture the whole face changes across frames in the temporal dimension. The decoder remaps the learned codes to the original face input size. The temporal consistency information for each landmark is used to improve alignment accuracy. It also consists of a two-layer RNN in between the encoder-decoder module. The first layer captures spatial-temporal appearance features whereas the second layer memorizes the temporal information across frames. Facial landmarks are determined by a weighted fusion of both spatial and temporal streams. The landmark positions are refined simultaneously in both the streams.

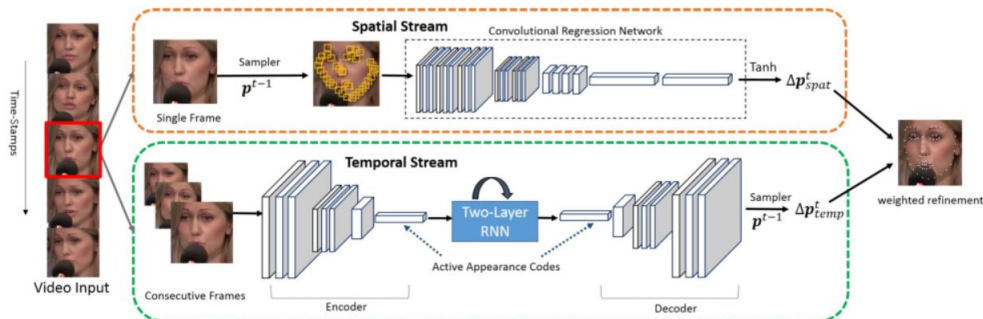


Fig. 6: TSTN pipeline. Source: [15]

124 3.5 Face Alignment Recurrent Network (2017) [10]

125 Previous state-of-the-art regression-based approaches start with an initial shape estimation and it-
 126 iteratively estimate the facial shape at successive stages by estimating an increment from the previous
 127 estimation. This paper proposes to improve the cascade shape regression by using LSTM and Re-
 128 gion Convolutional Neural Network (RCNN). The LSTM model exploits both spatial and temporal
 129 information for landmark detection in images and videos in uncontrolled conditions. The predicted
 130 landmark location is used as a basis for estimation in the next stage and frame in spatial and tem-
 131 poral dimensions. The process continues recurrently until the face shape is finalized. Figure 7 shows
 132 the training architecture of Face Alignment Recurrent Network. The face image, initial face shape,
 133 and ground truth shape are given as input to the network. The image is passed through several
 134 convolutional and max-pooling layers to obtain a feature map. The initial face shape contains facial
 135 landmarks. Region of Interest (ROI) pooling is applied around each landmark to obtain ROI pooling
 136 features. The ROI pooling features are concatenated and given to a fully connected layer followed
 137 by an LSTM layer. The network outputs the predicted shape increment over the initial face shape.
 138 The initial face shape is summed over the predicted shape increment to obtain updated initial face
 139 shape. This process continues recurrently for T stages.

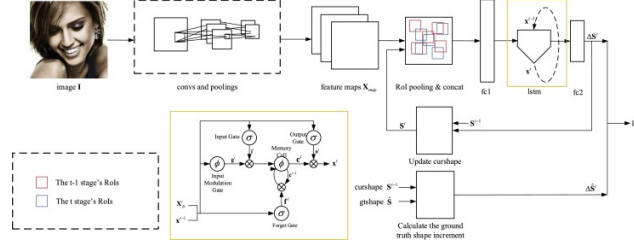


Fig. 7: FARN training architecture. Source: [10]

140 4 Performance Comparison

141 In this section we compare the above methods starting from the datasets used for training and
 142 testing, evaluation metrics, evaluation on common dataset and robustness to challenges.

143 **REDNet** [16] is trained on both image and video datasets with different configurations for
 144 different datasets. The training happens in 3 steps. In the first step, the network without the temporal
 145 recurrent learning and supervised identity disentangling modules is pre-trained using the image
 146 datasets AFLW, Helen and LFPW. In the second step, supervised identity disentangling is included
 147 and trained with other modules using the image-based LFW dataset. In the third step, the temporal
 148 recurrent learning module is included and the entire model is fine-tuned using the video dataset 300-
 149 VW. Inter-ocular distance is used to normalize the Root Mean Square Error(RMSE). The coarse-to-
 150 fine strategy in the two-step spatial recurrent learning(landmark detection and landmark regression)
 151 makes the model robust to large head pose and partial occlusions.

152 **Dynamic Facial Analysis** [8] is trained using the created SynHead dataset with the L2 loss
 153 function and tested on the BIWI dataset. It is then fine-tuned using training data from the BIWI
 154 dataset. For landmark detection, the corresponding model is trained and tested using a randomly
 155 split 300-VW dataset. For each frame, the mean Euclidean distance of the 68 landmarks normalized
 156 by the diagonal distance of the ground truth box is computed. The metrics used for evaluation are
 157 *area under the curve* which is the area under the cumulative error distribution curve, and *failure*

158 *rate* which is the percentage of images whose errors are larger than a given threshold. The proposed
 159 RNN-based architecture is more robust to large head poses and occlusions compared to per-frame
 160 estimation.

161 **DADRL** [9] is trained in two stages. The **first stage** is the *supervised learning stage* in which
 162 the two agents are trained separately. All training data from 300 faces in the wild challenge (300-W
 163 image dataset) is used to train the alignment agent. The 300-VW training set is used to train the
 164 tracking agent. The communicated messages are set to zero in this stage. The **second stage** is the
 165 *reinforcement learning stage* in which the whole network is trained with the 300-VW training set.
 166 The model is evaluated on the test set of 300-VW. For evaluation, Normalized Root Mean Square
 167 Error and cumulative error distribution plots are used. The communication channels between the
 168 agents has shown to provide robustness to occlusions. The authors also propose a DADRL-3D which
 169 is trained on the 3D Menpo dataset [20]. It is more robust to large pose as it is trained on 3D data.

170 **TSTN** [15] is trained using the 300-VW training set. The pre-trained spatial stream network
 171 is finetuned beforehand. The model is evaluated on the testing sets of Talking Face(TF) and 300-
 172 VW datasets. Normalized Root-Mean-Square-Error and cumulative error distribution plots are used
 173 for evaluating the model. The temporal stream consisting of the encoder-decoder module and the
 174 2-layer RNN provides consistency over time, and the spatial stream which provide complementary
 175 appearance information from still frames, achieve robustness to large head poses, occlusions and
 176 variations of expressions.

177 **FARN** [10] is trained on the training partition consisting of training sets of LFPW, Helen and the
 178 entire AFW with 3148 images in total. The testing partition contains 3 parts - the common subset,
 179 the challenging subset, and the full set. The common subset consists of testing set of LFPW and
 180 Helen with 554 images in total. The challenging subset consists of the IBUG dataset which contains
 181 additional annotations for 135 images in difficult poses and expressions. The full set consists of
 182 both the common subset and the challenging subset with 689 images. The model is evaluated using
 183 point-to-point Root Mean Square Error between the face shape and the ground truth annotations.
 184 The end-to-end trained model runs extremely fast (18ms) with robustness to large head poses and
 185 occlusions.

186 **REDNet** [16], **Dynamic Facial Analysis** [8], **TSTN** [15] and **FARN** [10] provide testing results on
 187 challenging category of 300-VW test set for 68 landmarks. **REDNet** [16] and **TSTN** [15](7 landmarks)
 188 provide results on Talking Face dataset [1]. **REDNet** [16] provides results for both 68 and 7 landmarks
 189 in both datasets. Table 2 reports the RMSE of the compared methods on 300-VW and TF [1]
 190 datasets. **DADRL** [9] uses normalized point-to-point error for evaluation and hence the evaluation
 191 results provided in the paper cannot be used to compare the method with other methods. **REDNet**
 192 [16] performs the best on the 300-VW dataset and has the least error among the compared methods
 193 for both 7 and 68 landmarks. **TSTN** [15] shows the best performance on the controlled TF dataset
 194 for 7 landmarks.

Table 2: Evaluation on 300-VW and TF test sets

Method	300-VW		TF	
	RMSE(68 landmarks)	RMSE(7 landmarks)	RMSE(68 landmarks)	RMSE(7 landmarks)
REDNet [16]	5.15	5.29	2.77	2.89
Dynamic Facial Analysis [8]	6.16			
DADRL [9]				
TSTN [15]	5.59			2.13
FARN [10]	5.90			

5 Conclusion

In this paper, we have reviewed some of the state-of-the-art deep learning methods for video-based face alignment. All of these methods avoid hand-engineering by using neural networks. Each of the methods is independent and not an improvisation of the other. All these methods use RNN in common to model temporal information. Although the recent methods have shown robustness to large head poses and occlusions, face tracking under difficult illumination is still a challenge.

References

1. Fgnet: Talking face video. *Tech. rep.*, 2004.
2. Christos Sagonas a, Epameinondas Antonakosa, Georgios Tzimiropoulosb, Stefanos Zafeirioua, and Maja Pantic. 300 faces in-the-wild challenge: database and results, 2016.
3. Jeremy Bailenson, Emmanuel (Manos) Pontikakis, Iris Mauss, James Gross, Maria Jabon, Cendri Hutcherson, Clifford Nass, and Oliver John. Real-time classification of evoked emotions using facial feature tracking and physiological responses. *International Journal of Human-Computer Studies*, 66:303–317, 05 2008.
4. Peter Belhumeur, David Jacobs, David Kriegman, and Neeraj Kumar. Localizing parts of faces using a consensus of exemplars. *IEEE transactions on pattern analysis and machine intelligence*, 35:2930–40, 12 2013.
5. X. P. Burgos-Artizzu, P. Perona, and P. Dollár. Robust face landmark estimation under occlusion. In *2013 IEEE International Conference on Computer Vision*, pages 1513–1520, Dec 2013.
6. Paola Campadelli, Raffaella Lanzarotti, and C. Savazzi. A feature-based face recognition system. pages 68–73, 10 2003.
7. Gabriele Fanelli, Matthias Dantone, Juergen Gall, Andrea Fossati, and Luc Van Gool. Random forests for real time 3d face analysis. *International Journal of Computer Vision*, 101, 02 2013.
8. Jinwei Gu, Xiaodong Yang, Shalini Mello, and Jan Kautz. Dynamic facial analysis: From bayesian filtering to recurrent neural network. pages 1531–1540, 07 2017.
9. Minghao Guo, Jiwen Lu, and Jie Zhou. *Dual-Agent Deep Reinforcement Learning for Deformable Face Tracking: 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part X*, pages 783–799. 09 2018.
10. Qiqi Hou, Jinjun Wang, Ruibin Bai, Sanping Zhou, and Yihong Gong. Face alignment recurrent network. *Pattern Recognition*, 74, 09 2017.
11. Gary Huang, Marwan Mattar, Tamara Berg, and Eric Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. *Tech. rep.*, 10 2008.
12. Martin Köstinger, Paul Wohlhart, Peter M. Roth, and Horst Bischof. Annotated facial landmarks in the wild: A large-scale, real-world database for facial landmark localization. pages 2144–2151, 11 2011.
13. Vuong Le, Jonathan Brandt, Zhe Lin, Lubomir Bourdev, and Thomas Huang. Interactive facial feature localization. 10 2012.
14. Feng Liu, Qijun Zhao, xiaoming Liu, and Dan Zeng. Joint face alignment and 3d face reconstruction with application to face recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, page 1–1, 2018.
15. Hao Liu, Jiwen Lu, Jianjiang Feng, and Jie Zhou. Two-stream transformer networks for video-based face alignment. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PP:1–1, 08 2017.
16. Xi Peng, Rogerio S. Feris, Xiaoyu Wang, and Dimitris N. Metaxas. Red-net: A recurrent encoder-decoder network for video-based face alignment, 2018.
17. Xi Peng, Shaoting Zhang, Yang Yu, and Dimitris Metaxas. Piefa: Personalized incremental and ensemble face alignment. 12 2015.
18. Jie Shen, Stefanos Zafeiriou, Grigorios Chrysos, Jean Kossaifi, Georgios Tzimiropoulos, and Maja Pantic. The first facial landmark tracking in-the-wild challenge: Benchmark and results. 12 2015.
19. Keqiang Sun, Wayne Wu, Tinghao Liu, Shuo Yang, Quan Wang, Qiang Zhou, Zuochang Ye, and Chen Qian. Fab: A robust facial landmark detection framework for motion-blurred videos, 2019.
20. Stefanos Zafeiriou, Grigorios G Chrysos, Anastasios Roussos, Evangelos Ververas, Jiankang Deng, and George Trigeorgis. The 3d menpo facial landmark tracking challenge. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 2503–2511, 2017.