



Unsupervised Performance Analysis of 3D Face Alignment

Mostafa Sadeghi, Sylvain Guy, Adrien Raison, Xavier Alameda-Pineda, Radu Horaud

► To cite this version:

Mostafa Sadeghi, Sylvain Guy, Adrien Raison, Xavier Alameda-Pineda, Radu Horaud. Unsupervised Performance Analysis of 3D Face Alignment. 2020. <hal-02543069v3>

HAL Id: hal-02543069

<https://hal.science/hal-02543069v3>

Preprint submitted on 16 Oct 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

Unsupervised performance analysis of 3D face alignment

Mostafa **Sadeghi**^a, Sylvain **Guy**^b, Adrien **Raison**^a, Xavier **Alameda-Pineda**^a, Radu **Horaud**^{a,**}

^a*Inria Grenoble, Montbonnot Saint-Martin, FRANCE*

^b*Université Grenoble Alpes, FRANCE*

ABSTRACT

We address the problem of analyzing the performance of 3D face alignment (3DFA) algorithms. Traditionally, performance analysis relies on carefully annotated datasets. Here, these annotations correspond to the 3D coordinates of a set of pre-defined facial landmarks. However, this annotation process, be it manual or automatic, is rarely error-free, which strongly biases the analysis. In contrast, we propose a fully unsupervised methodology based on robust statistics and a parametric confidence test. We revisit the problem of robust estimation of the rigid transformation between two point sets and we describe two algorithms, one based on a mixture between a Gaussian and a uniform distribution, and another one based on the generalized Student's t-distribution. We show that these methods are robust to up to 50% outliers, which makes them suitable for mapping a face, from an unknown pose to a frontal pose, in the presence of facial expressions and occlusions. Using these methods in conjunction with large datasets of face images, we build a statistical frontal facial model and an associated parametric confidence metric, eventually used for performance analysis. We empirically show that the proposed pipeline is neither method-biased nor data-biased, and that it can be used to assess both the performance of 3DFA algorithms and the accuracy of annotations of face datasets.

1. Introduction

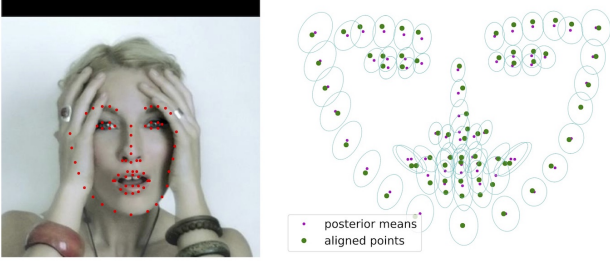
The problem of face alignment (FA) is the problem of facial landmark detection and localization from a single RGB image. Face alignment is an important research topic as it provides input to a variety of computer vision tasks, such as head-pose estimation and tracking, face recognition, facial expression understanding, visual speech recognition, etc., (Escalera et al, 2018; Loy et al, 2019). 2D face alignment (2DFA) has been extensively studied for the last decades, yielding a plethora of methods and algorithms (Wu and Ji, 2019). State of the art 2DFA based on deep neural networks (DNNs) are the best-performing methods in terms of accuracy, invariance with respect to facial appearances, shapes, expressions, as well as in terms of repeatability and reproducibility in the presence of image noise, image resolution, motion blur, lighting conditions and varying backgrounds.

Nevertheless, 2DFA methods yield poor landmark detection and localization performance in the presence of occlusions

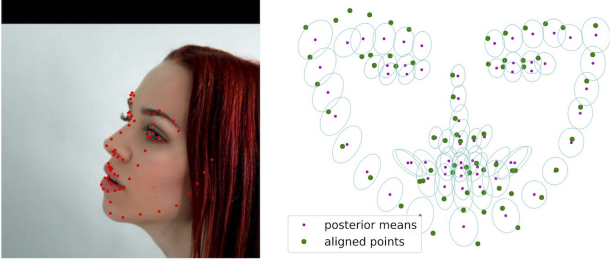
which occur in case of large poses induced by out-of-image-plane head rotations (self occlusions) as well as by the presence of various objects in the camera field of view, such as glasses, hair, hands and handheld objects, etc. Robust facial landmark detection and localization in the presence of occlusions can only be achieved on the premise that 3D information is taken into account. It is well established that 2D facial landmarks (and, more generally, face images) embed 3D information. This information can be retrieved by fitting a 3D face model to a 2D face image, even if the latter is only partially visible. The process of fitting a 3D model to a 2D image constitutes the basis of training 3D face alignment (3DFA) algorithms.

Consider for example a 3D face model that is parameterized both by identities and by facial deformations, e.g. the parametric 3D deformable model (3DMM) (Banz and Vetter, 1999). The task of fitting 3DMM to an RGB image of a face consists of estimating the parameters of the mapping from the 3D generic model to a particular face, namely the identity and expression parameters, as well as the pose parameters (scale, rotation, translation and projection), e.g. (Gou et al, 2016; Zhu et al, 2016). Once an optimal set of parameters is found, one can associate 3DMM vertices with facial landmarks. This stays

^{**}Corresponding author: Tel.: +33 476 615226;
e-mail: Radu.Horaud@inria.fr (Radu Horaud)



(a) Analysis of 3D landmarks extracted with (Bulat and Tzimiropoulos, 2016)



(b) Analysis of ground-truth 3D landmarks from the AFLW2000-3D dataset (Zhu et al, 2016)

Fig. 1: Two examples from the AFLW2000-3D dataset (Zhu et al, 2016) (left). The landmarks are mapped onto a statistical frontal landmark model (right) built using the YawDD dataset (Abtahi et al, 2014) and (Feng et al, 2018), which enables us to verify whether the mapped landmarks lie within their associated ellipsoidal confidence volumes or not.

at the basis of many automatic and semi-automatic methods for annotating 2D faces with 3D landmarks, e.g. (Deng et al, 2019).

Nevertheless, the fitting task just mentioned is a difficult non-linear optimization problem, in particular in the presence of large poses and of occlusions. In the recent past, a number of methods has been developed to perform this 3D-to-2D fitting process necessary for 3D facial landmark annotation. The performance of the vast majority of existing 3DFA methods inherently rely on the quality of landmark annotation. This is true for training using modern discriminative deep learning methods, but it is true for testing as well. Indeed, to date, algorithm performance is computed empirically by measuring the error between the predicted output and the corresponding ground-truth, e.g. (Jeni et al, 2016; Deng et al, 2019). Under these circumstances, annotation errors are likely to bias both parameter estimation (training) and performance evaluation (testing).

There is a lack of a benchmarking methodology that could assess *quantitatively* and in a completely unsupervised manner the robustness and effectiveness of 3DFA algorithms, namely a method that computes a confidence score that measures algorithm performance in the absence of the ground truth. This is also crucial in order to decide *without human intervention*, whether a 3DFA method, that is applied to an unknown image of a face with no annotation available, yields an output that is accurate enough to be further used by other algorithms, such as head-gaze estimation, facial expression analysis or lip reading.

This paper proposes a methodological framework for assessing the performance of 3DFA algorithms based on *robust probability distribution functions* and on a *statistical confidence test*.

Unlike supervised metrics currently in use for 3DFA performance evaluation and based on annotated datasets, the proposed method is fully unsupervised. We show that the robust estimation of the rigid mapping between two sets of 3D facial landmarks, namely (i) a predicted set, associated with a face with unknown identity, pose and expression, and (ii) a model set associated with a statistical frontal face, provides a reliable way to separate face pose (due to head motions) from non-rigid face deformations (due to facial expressions), all in the presence of badly located landmarks.

Using a 3DFA algorithm and a very large and unannotated dataset of face images with large variabilities in pose, expression and identity, we make use of the robust rigid-mapping methodology to build a *statistical frontal landmark model* and a parameterized *confidence score*. Based on this pipeline, the proposed performance evaluation protocol proceeds as follows. First, 3D landmarks are extracted from a face image using a 3DFA method. Second, the landmarks are rigidly mapped onto the frontal model. Third, a confidence score is computed for each mapped landmark, thus allowing to assess whether the landmark lies within a *confidence volume* or not.

We describe in detail an experimental evaluation framework that uses publicly available datasets and 3DFA software packages associated with three published articles and one unpublished paper. We empirically show that our methodological pipeline is neither dataset- nor method-biased. We also show that the proposed framework can be used not only to assess quantitatively the performance of 3DFA algorithms, but also to test the accuracy of automatic and semi-automatic methods currently used for the annotation of face datasets.

The methodology proposed in this paper is illustrated in Fig. 1. The two images (left) are from the AFLW2000-3D dataset (Zhu et al, 2016). The statistical frontal landmark model (right) is built using the 3DFA method of (Feng et al, 2018) and the YawDD dataset (Abtahi et al, 2014). This model characterizes each landmark with an ellipsoidal confidence volume centered at a posterior mean. Fig. 1(a): Landmarks extracted using (Bulat and Tzimiropoulos, 2016) (left) and mapped onto the statistical model (right). In this case, most of the landmarks lie inside their confidence volumes, thus assessing their correctness. Fig. 1(b): Ground-truth landmarks obtained with a semi-automatic annotation process (Zhu et al, 2016) and mapped onto the statistical model (right). One may notice that in this case, many mapped landmarks fall outside their confidence volumes. The benefit of the proposed method is twofold: (i) an unsupervised assessment of the quality of the detected landmarks, and (ii) a robust and expression- and identity-preserving landmark mapping from an arbitrary pose to a frontal pose.

The remainder of this paper is organized as follows. Section 2 reviews the related work. Section 3 summarizes the problem formulation. Section 4 briefly reviews maximum likelihood estimation (MLE) for rigid mapping and describes two robust methods. Section 5 empirically analyses the proposed rigid-mapping methods. Section 6 proposes a methodological pipeline for building a statistical face model and an associated parametric confidence metric. Section 7 presents extensive ex-

perimental results, and Section 8 draws some conclusions.¹

2. Related Work

It is interesting to note that the recently proposed methods for 3DFA lie at the crossroads of deformable shape models, model-based image analysis and neural networks. In order to discuss these links, we introduce some mathematical notations and concepts. Let vector $\mathbf{p} \in \mathcal{P} \subset \mathbb{R}^K$ denote the ensemble of parameters of a 3D face model S (identity, expression and pose), where \mathcal{P} is the parameter vector space, K is the number of parameters, and let $\mathbf{b} \in \mathcal{B} \subset \mathbb{R}^{I \times J}$ denote the image of a face from a set of images of size $I \times J$. One class of 3DFA methods directly learns a mapping $\mathbf{p} = f(\mathbf{b})$ from a training dataset of M face images and their associated model parameters $\{\mathbf{b}_m, \mathbf{p}_m\}_{m=1}^M$, e.g. (Zhu et al, 2016; Jourabloo and Liu, 2017; Feng et al, 2018; Deng et al, 2019).

Another class of methods proceeds in two steps. First, 2D landmarks are extracted from a face image by learning an image-to-landmark mapping $\mathbf{u} = g(\mathbf{b})$, from a face image to a set of 2D landmarks $\mathbf{u} = \{\mathbf{v}_n\}_{n=1}^N \in \mathbb{N}^{2 \times N}$, and using a training dataset $\{\mathbf{b}_m, \mathbf{u}_m\}_{m=1}^M$. Second, a 2D-to-3D mapping $\mathbf{s} = h(\mathbf{u})$ is estimated, where $\mathbf{s} = \{\mathbf{x}_n\}_{n=1}^N \in \mathbb{R}^{3 \times N}$ is a set of N 3D landmarks. This mapping can be obtained either by learning, using a training dataset conditioned by a 3D model parameterized by \mathbf{p} , i.e. $\{\mathbf{u}_m, \mathbf{s}(\mathbf{p}_m)\}_{m=1}^M$, e.g. (Zhao et al, 2016; Bulat and Tzimiropoulos, 2016, 2017), or by direct optimization over \mathbf{p} of a function that maps a 3D model onto the 2D landmarks, e.g. (Yu et al, 2017).

These 3DFA DNN-based methods use a variety of architectures in order to learn the regression functions $f(\cdot)$, $g(\cdot)$ and $h(\cdot)$ mentioned above. Given this variety, it is difficult to directly compare them and assess their merits based on the analysis of the underlying DNN concepts and methodologies. Hence, 3DFA algorithm performance should be measured empirically, as is often the case in deep learning.

To date and to the best of our knowledge, there has been a handful of 3DFA benchmark datasets and corresponding evaluation metrics, (Jeni et al, 2016; Deng et al, 2019; Sanyal et al, 2019). In (Jeni et al, 2016), four datasets were specifically gathered, annotated and prepared, and two performance metrics were used for this challenge. The BU-4DFE (Yin et al, 2008) and BP-4D-Spontaneous (Zhang et al, 2014) datasets used a structured-light stereo sensor to capture textured 3D meshes of faces in controlled conditions and with various backgrounds. 2295 meshes were selected from these datasets and manually annotated with 66 landmarks and with self-occlusion information. Then, 16065 2D views were synthesized (seven views for each mesh) with yaw and pitch rotations ranging in the intervals $[-45^\circ, +45^\circ]$ and $[-30^\circ, +30^\circ]$, respectively. Additionally, there were 7,000 frames from the Multi-PIE (Gross et al,

2010) and 541 frames from the Time-Sliced datasets, respectively. Both these datasets contain RGB images gathered with multiple cameras from different viewpoints but with no 3D information. Therefore, a 3D face model is extracted for each image, using the model-based multi-view structure-from-motion technique of (Jeni et al, 2017). As above, each 3D face model was annotated with 66 landmarks and with self-occlusion information.

The following metrics were used in (Jeni et al, 2016): the ground-truth error (GTE) and the cross-view ground-truth consistency error (CVGTCE), namely,

$$\text{GTE}(i) = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_{n,i} - \hat{\mathbf{x}}_{n,i}\|/d_i, \quad (1)$$

$$\text{CVGTCE}(i) = \frac{1}{N} \sum_{n=1}^N \|s_i \mathbf{R}_i \mathbf{x}_{n,i} + \mathbf{t}_i - \hat{\mathbf{x}}_{n,i}\|/d_i, \quad (2)$$

where $\mathbf{x}_{n,i}$ denotes the n -th detected 3D landmark associated with test sample i , $\hat{\mathbf{x}}_{n,i}$ is the corresponding ground-truth 3D landmark, d_i is the inter-ocular distance of the sample face i , N is the number of landmarks, and s_i , \mathbf{R}_i , and \mathbf{t}_i are the scale factor, rotation matrix and translation vector associated with a rigid mapping that compensates the possible discrepancy between the set of detected landmarks and the set of ground-truth landmarks.

The Menpo challenge (Deng et al, 2019) is based on a dataset of over 12000 face images. In order to obtain 2D and 3D ground-truth landmarks, an automatic annotation process is proposed, which fits a 3D face model to each 2D image. This fitting is carried out via non-linear minimization over the shape parameters (identity and expression), the rigid parameters (rotation and translation of the 3D model with respect to the camera), and the camera parameter (the scale of the weak-perspective model). The evaluation metric uses (1) with a different normalization factor, namely the size of the face bounding box.

The NoW benchmark is proposed in (Sanyal et al, 2019) for the task of 3D reconstruction from a single monocular image of a face. The associated dataset contains 2054 face images in frontal and profile views of 100 subjects and a 3D head scan for each subject. This dataset is similar in spirit to (Bagdanov et al, 2011). While the images contain four categories (neutral, expression, occlusion, and selfie) the 3D scans correspond to neutral faces. Therefore, the challenge for all categories is the reconstruction a neutral 3D face, which implies that non-neutral faces must undergo some form of disentanglement. Moreover, since the predicted 3D mesh and the ground-truth 3D scan lie in different coordinate systems, a rigid alignment must be performed prior to computing an evaluation metric. The authors provide an alignment procedure that minimizes a scan-to-mesh (or point-to-surface) distance over the alignment parameters (scale, rotation, and translation). This is a difficult alignment problem that necessitates to alternate between (i) selecting the closest points on the mesh and (ii) estimating the rigid parameters. Once an optimal alignment is found, the evaluation metric consists of the scan-to-mesh distance.

¹Supplemental material for this paper can be found at <https://team.inria.fr/perception/research/upa3dfa/>

The evaluation metrics used in these benchmarks require ground-truth either of 3D landmarks or of 3D scans. Manual annotations of thousands of images cannot be achieved and automatic annotation must therefore be used. As outlined above, automatic annotation is based on complex non-linear minimization methods that are prone to errors and may not be reliable in the presence of profile views and of occlusions. Localization noise is inherent. Nevertheless, these evaluation metrics are limited in scope since they cannot distinguish between landmark localization noise (inlying data) and large localization errors (outlying data).

In contrast, the proposed methodology does not make use of ground-truth annotations. Robust rigid alignment (analyzed in detail below) is used to build frontal landmark models in a completely unsupervised way. A statistical characterization of each landmark is provided by measuring the discrepancy between the predicted landmark and the corresponding model landmark. Indeed, a confidence score is computed for each predicted landmark in order to assess its localization accuracy and to decide whether the landmark should be treated as an inlier (affected by detection noise) or as an outlier (the detection has failed). Our method may well be viewed as an analysis of performance of 3DFA algorithms, rather than a benchmark or a challenge. This is particularly useful whenever the output of 3DFA is used for facial expression recognition, for lip reading, for 3D pose estimation, etc. The proposed landmark analysis can also be applied to ground-truth landmarks in order to remove bad annotations, be them manual, semi-automatic or fully automatic.

A fundamental building block of the proposed methodology is the estimation of the rigid transformation (scale, rotation and translation) between two 3D point sets. We propose to perform this estimation in a robust way, where the term robustness refers to the capacity of a method to be unaffected by outliers. For that purpose, we cast the problem at hand in the framework of MLE, i.e. (7). In MLE, the choice of the likelihood function is crucial. We opt for two choices, namely a mixture model with two components, a Gaussian component and a uniform one (GUM) (Banfield and Raftery, 1993; Zaharescu and Horaud, 2009; Lathuilière et al, 2018) and the generalized Student’s t-distribution (McLachlan and Peel, 2000; Sun et al, 2010; Forbes and Wraith, 2014). These two distributions behave quite differently. For each data point, GUM evaluates the posterior probability of being an inlier or an outlier, i.e. (16) and (18). Student evaluates a weight w associated with each data point. These weights play the role of precisions (the inverse of the variance): the higher the weights, the more reliable data. Each weight is treated as a random variable modeled with a gamma distribution, i.e. (24). Generalized Student belongs to the larger class of heavy-tailed distributions that are well-known to be robust to outliers (McLachlan and Peel, 2000). To the best of our knowledge, this is the first time that a heavy-tailed distribution is used in conjunction with rigid alignment.

Both GUM and Student evaluate a posterior covariance matrix, i.e. (23) and (29), respectively. The evaluation of these covariance matrices is fundamental for taking into account heterogeneous landmark distributions, e.g. Fig. 1-right, and for as-

sessing how much one can trust the results. This stays in strong contrast with the prevailing methods that are used in computer vision to estimate the rigid transformation between two point sets (Horn, 1987; Horn et al, 1988; Faugeras and Hebert, 1986; Arun et al, 1987; Umeyama, 1991). These methods assume an isotropic covariance matrix, i.e. spherical uncertainties, thus leading to closed-form solutions for estimating the rigid parameters. The fact that both GUM and Student estimate full covariance matrices, i.e. ellipsoidal uncertainties, introduces an additional complexity, namely a non-linear solver is required to estimate the rotation matrix (please consult Section 4.1). In (Horaud et al, 2010) this was addressed via convex relaxation. In this paper, we simplify the optimization problem by using quaternions and propose to use sequential quadratic programming (Bonnans et al, 2006) (please refer to Section 4.4 for an in-depth discussion).

The proposed confidence test detects anomalous outputs predicted by a 3DFA algorithm, be it based on a deep architecture or not. Alternatively, a number of methods were recently proposed to detect anomalous inputs, most notably out-of-distribution (OOD) detection methods, e.g. (Hendrycks and Gimpel, 2017; Liang et al, 2018; Lee et al, 2018; Hendrycks et al, 2019). These methods are concerned with classification networks that use a softmax layer to predict a probability distribution over the discrete set of possible labels (classes). Moreover these methods require annotated OOD examples. In contrast our proposed confidence test applies to regression networks and is unsupervised.

3. Problem Formulation

Let us consider the mapping between two sets of 3D facial landmarks, a predicted set, $\mathbf{x}_{1:N} = (\mathbf{x}_1 \dots \mathbf{x}_N) \in \mathbb{R}^{3 \times N}$, and a model set, $\mathbf{y}_{1:N} = (\mathbf{y}_1 \dots \mathbf{y}_N) \in \mathbb{R}^{3 \times N}$. The predicted set corresponds to a face with arbitrary and unknown pose, identity, expression and occlusion. Without loss of generality, the model set corresponds to a *statistical frontal landmark model*, namely each landmark n in this set is characterized by a mean vector \mathbf{y}_n and a covariance matrix $\mathbf{C}_n \in \mathbb{R}^{3 \times 3}$, i.e. Section 6.2. In the general case, the unknown-to-frontal mapping is composed of a rigid transformation, i.e. head motion, and of a non-rigid facial deformation, up to an unknown error. Therefore, we introduce an additive residual $\mathbf{r}_n \in \mathbb{R}^3$, associated with the n th landmark, to account for the non-rigid component of the mapping and for various sources of errors. The mapping, from an unknown face pose to a frontal face pose can be modeled in the following way:

$$\mathbf{y}_n = s\mathbf{R}\mathbf{x}_n + \mathbf{t} + \mathbf{r}_n, \forall n \in \{1, \dots, N\}, \quad (3)$$

where the rigid transformation is parameterized by a scale factor $s \in \mathbb{R}$, a rotation matrix $\mathbf{R} \in SO(3) \subset \mathbb{R}^{3 \times 3}$ and a translation vector $\mathbf{t} \in \mathbb{R}^3$, while the non-rigid deformations and errors are represented by the residuals $\mathbf{r}_{1:N}$. Let us assume that the optimal scale, rotation and translation are obtained via MLE, or equivalently, minimization of the negative log-likelihood pa-

parameterized by θ :

$$\{s^*, \mathbf{R}^*, \mathbf{t}^*, \theta^*\} = \underset{s, \mathbf{R}, \mathbf{t}, \theta}{\operatorname{argmin}} \sum_{n=1}^N -\log P(\mathbf{r}_n; \theta), \quad (4)$$

followed by an estimation of the residuals:

$$\mathbf{r}_n^* = \mathbf{y}_n - s^* \mathbf{R}^* \mathbf{x}_n - \mathbf{t}^*, \forall n \in \{1, \dots, N\}. \quad (5)$$

Note that the minimizer (4) must be immune to non-rigid deformations and to noise, which are jointly referred to as *inliers*, and at the same time it must be robust to large errors in landmark localization due to various perturbations, such as occlusions, motion blur, etc., which are referred to as *outliers*, i.e. Section 4. Therefore, we seek a *robust rigid mapping* technique that enables us to discriminate between inliers and outliers, namely:

$$\begin{cases} \text{if } \|\mathbf{r}_n^*\|_{\widetilde{\mathbf{C}}_n} \leq 1 & n = \text{inlier} \\ \text{otherwise} & n = \text{outlier}, \end{cases} \quad (6)$$

where $\widetilde{\mathbf{C}}_n = 9\mathbf{C}_n$ and $\|\cdot\|_{\mathbf{C}}$ is the Mahalanobis distance associated with covariance \mathbf{C} . This guarantees with 99% confidence that if the n th landmark lies inside the ellipsoid defined by $\widetilde{\mathbf{C}}_n$, then it is an inlier, and otherwise it is an outlier (it lies outside the ellipsoid), i.e. Section 6.3. The aptitude to discriminate between inlying and outlying landmarks is a core feature of the proposed unsupervised performance analysis.

4. Robust Rigid Mapping

We cast the problem of estimating the rigid mapping between two sets of landmarks into the framework of robust probability distribution functions. We assume that the residuals $\mathbf{r}_{1:N}$ are independent and identically distributed (i.i.d). Then, the problem of estimating the rigid transformation parameters can be solved via MLE or, equivalently, via negative log-likelihood minimization, namely

$$\min_{\theta} \mathcal{L}(\theta | (\mathbf{x}, \mathbf{y})_{1:N}),$$

with:

$$\mathcal{L}(\theta | (\mathbf{x}, \mathbf{y})_{1:N}) = -\sum_{n=1}^N \log P(\mathbf{r}_n; \theta), \quad (7)$$

where $P(\mathbf{r}; \theta)$ is the probability distribution function (pdf) of \mathbf{r} parameterized by θ which is composed of s , \mathbf{R} , \mathbf{t} and of the pdf parameters that are specified below.

4.1. Gaussian Model

The simplest statistical model is to assume that the residuals follow a zero-centered Gaussian distribution with covariance matrix $\Sigma \in \mathbb{R}^{3 \times 3}$, namely $P(\mathbf{r}; \theta) = \mathcal{N}(\mathbf{r}; \mathbf{0}, \Sigma)$. By developing (7) and ignoring terms that do not depend on the model parameters, we obtain:

$$\mathcal{L}(\theta | (\mathbf{x}, \mathbf{y})_{1:N}) = \frac{1}{2} \sum_{n=1}^N (\|\mathbf{y}_n - s\mathbf{R}\mathbf{x}_n - \mathbf{t}\|_{\Sigma}^2 + \log |\Sigma|), \quad (8)$$

where $\|\mathbf{a}\|_{\Sigma}^2 = \mathbf{a}^T \Sigma^{-1} \mathbf{a}$ is the squared Mahalanobis norm of $\mathbf{a} \in \mathbb{R}^3$. The minimization of (8) over \mathbf{t} yields:

$$\mathbf{t}^* = \bar{\mathbf{y}} - s^* \mathbf{R}^* \bar{\mathbf{x}}, \quad (9)$$

where the over-script $*$ indicates the optimal value of a parameter and with the notations:

$$\bar{\mathbf{x}} = \frac{\sum_{n=1}^N \mathbf{x}_n}{N}, \quad \bar{\mathbf{y}} = \frac{\sum_{n=1}^N \mathbf{y}_n}{N}. \quad (10)$$

By substituting (9) into (8) and by using centered coordinates, i.e. $\mathbf{x}'_n = \mathbf{x}_n - \bar{\mathbf{x}}$, $\mathbf{y}'_n = \mathbf{y}_n - \bar{\mathbf{y}}$, we obtain:

$$\mathcal{L}(\theta | (\mathbf{x}', \mathbf{y}')_{1:N}) = \frac{1}{2} \sum_{n=1}^N (\|\mathbf{y}'_n - s\mathbf{R}\mathbf{x}'_n\|_{\Sigma}^2 + \log |\Sigma|). \quad (11)$$

Standard approaches to the minimization of (11) with respect to the rotation matrix assume an isotropic covariance, $\Sigma = \sigma^2 \mathbf{I}_3$. Indeed, the development of (11) yields

$$\sum_{n=1}^N s\mathbf{x}'_n \mathbf{R} \Sigma^{-1} \mathbf{R}^T \mathbf{x}'_n{}^T = s\sigma^{-2} \sum_{n=1}^N \mathbf{x}'_n \mathbf{x}'_n{}^T$$

thus leading to closed-form solutions, e.g. (Horn, 1987; Horn et al, 1988; Faugeras and Hebert, 1986; Arun et al, 1987; Umeyama, 1991) and Appendix Appendix A. Nevertheless, the isotropic-covariance assumption is barely valid in practice. In the case of a full covariance, the optimization becomes

$$\mathbf{R}^* = \underset{\mathbf{R}}{\operatorname{argmin}} \frac{1}{2} \operatorname{tr}(\Sigma^{-1} (s^2 \mathbf{R} \mathbf{A} \mathbf{R}^T - 2s\mathbf{R}\mathbf{B})), \quad (12)$$

where $\operatorname{tr}(\cdot)$ is the trace operator and with the notations $\mathbf{A} = \sum_{n=1}^N \mathbf{x}'_n \mathbf{x}'_n{}^T$, $\mathbf{B} = \sum_{n=1}^N \mathbf{x}'_n \mathbf{y}'_n{}^T$. A rotation matrix must satisfy $\mathbf{R}\mathbf{R}^T = \mathbf{I}_3$ and $|\mathbf{R}| = +1$. This yields a constrained non-linear optimization problem. An elegant formulation consists of parameterizing the rotation with a unit quaternion, thus reducing the number of parameters from nine to four, while the number of constraints is reduced from seven to one. Let $\mathbf{R}(\mathbf{q})$, where \mathbf{q} is a unit quaternion (please consult Appendix Appendix A). Using this representation, the rotation is described by four parameters and the associated constrained optimization problem writes:

$$\mathbf{q}^* = \underset{\mathbf{q}}{\operatorname{argmin}} \frac{1}{2} \left(\operatorname{tr}(\Sigma^{-1} (s^2 \mathbf{R}(\mathbf{q}) \mathbf{A} \mathbf{R}(\mathbf{q})^T - 2s\mathbf{R}(\mathbf{q})\mathbf{B})) + \lambda(\mathbf{q}^T \mathbf{q} - 1)^2 \right). \quad (13)$$

Similar to (Horn, 1987) the optimal scale factor has a closed-form expression:

$$s^* = \left(\frac{\sum_{n=1}^N \mathbf{y}'_n{}^T \Sigma^{-1} \mathbf{y}'_n}{\sum_{n=1}^N (\mathbf{R}^* \mathbf{x}'_n)^T \Sigma^{-1} \mathbf{R}^* \mathbf{x}'_n} \right)^{1/2}. \quad (14)$$

Finally, the optimal covariance is estimated with:

$$\Sigma^* = \frac{1}{N} \sum_{n=1}^N (\mathbf{y}'_n - s^* \mathbf{R}^* \mathbf{x}'_n)(\mathbf{y}'_n - s^* \mathbf{R}^* \mathbf{x}'_n)^T. \quad (15)$$

Once the rotation and scale are initialized using the method of (Horn, 1987), alternating optimization can be used by iterating between (13), (14) and (15). We will refer to this method as *generalized Horn*.

4.2. Gaussian-uniform Mixture Model

Unfortunately, the above statistical model does not behave well in the presence of large residuals, or outliers. For the purpose of explicitly modeling inliers and outliers, a discrete random variable Z_n is associated with each residual \mathbf{r}_n , and let z be a realization of Z . Now, \mathbf{r} is drawn either from a zero-centered Gaussian distribution, as above, or from a multivariate uniform distribution:

$$P(\mathbf{r}|Z = z) = \begin{cases} \mathcal{N}(\mathbf{r}; \mathbf{0}, \Sigma) & \text{if } z = \text{inlier} \\ \mathcal{U}(\mathbf{r}; 0, \gamma) & \text{if } z = \text{outlier}, \end{cases} \quad (16)$$

where γ is the volume of the distribution. This yields a two-component mixture model, an inlier component with prior probability p , and an outlier component with prior probability $1 - p$. This naturally leads to solving the problem via expectation-maximization (EM) which alternates between (i) evaluating the posterior probabilities of the residuals to be inliers or outliers, and (ii) minimizing the *expected complete-data negative log-likelihood*, $E_Z[-\log P(\mathbf{r}_{1:N}, Z_{1:N}|\mathbf{r}_{1:N}; \theta)]$, where the expectation is taken over the realizations of Z , and where the parameter vector is $\theta = \{s, \mathbf{R}, p, \Sigma\}$.² This yields the minimization of:

$$\mathcal{E}(\theta|(\mathbf{x}', \mathbf{y}')_{1:N}) = \frac{1}{2} \sum_{n=1}^N \alpha_n (\|\mathbf{y}'_n - s\mathbf{R}\mathbf{x}'_n\|_{\Sigma}^2 + \log |\Sigma|) \quad (17)$$

where the posterior probability to be an inlier $\alpha_n = P(Z = \text{inlier}|\mathbf{r}_n)$, is :

$$\alpha_n = \frac{p \mathcal{N}(\mathbf{r}_n; \mathbf{0}, \Sigma)}{p \mathcal{N}(\mathbf{r}_n; \mathbf{0}, \Sigma) + (1 - p) \gamma^{-1}}, \quad (18)$$

and the posterior probability to be an outlier is $P(Z = \text{outlier}|\mathbf{r}_n) = 1 - \alpha_n$. The presence of $\alpha_{1:N}$ in (17) replaces (10) with:

$$\bar{\mathbf{x}} = \frac{\sum_{n=1}^N \alpha_n \mathbf{x}_n}{\sum_{n=1}^N \alpha_n}, \quad \bar{\mathbf{y}} = \frac{\sum_{n=1}^N \alpha_n \mathbf{y}_n}{\sum_{n=1}^N \alpha_n}, \quad (19)$$

as well as \mathbf{A} and \mathbf{B} from (13) with

$$\mathbf{A} = \sum_{n=1}^N \alpha_n \mathbf{x}'_n \mathbf{x}'_n{}^\top, \quad \mathbf{B} = \sum_{n=1}^N \alpha_n \mathbf{x}'_n \mathbf{y}'_n{}^\top. \quad (20)$$

Hence, (13) can be used to estimate the optimal rotation. Moreover, (14) is replaced with :

$$s^* = \left(\frac{\sum_{n=1}^N \alpha_n \mathbf{y}'_n{}^\top \Sigma^{-1} \mathbf{y}'_n}{\sum_{n=1}^N \alpha_n (\mathbf{R}^* \mathbf{x}'_n)^\top \Sigma^{-1} \mathbf{R}^* \mathbf{x}'_n} \right)^{1/2}. \quad (21)$$

The prior probability p and covariance matrix Σ are estimated with:

$$p = \frac{1}{N} \sum_{n=1}^N \alpha_n, \quad (22)$$

$$\Sigma^* = \frac{\sum_{n=1}^N \alpha_n (\mathbf{y}'_n - s^* \mathbf{R}^* \mathbf{x}'_n)(\mathbf{y}'_n - s^* \mathbf{R}^* \mathbf{x}'_n)^\top}{\sum_{n=1}^N \alpha_n}. \quad (23)$$

We refer to this model as the *Gaussian-uniform mixture* (GUM) and the associated EM is summarized in Algorithm 1.

Data: Centered point coordinates, i.e. (10).

Normalization parameter γ ;

Initialization of θ^{old} : Use the closed-form solution (Horn, 1987) to evaluate s^{old} and \mathbf{R}^{old} and then use these parameter values to evaluate Σ^{old} using (15) and set $p^{\text{old}} = 0.8$;

while $\|\theta^{\text{new}} - \theta^{\text{old}}\| > \epsilon$ **do**

E-step: Evaluate the posteriors $\alpha_{1:N}$ using (18) with θ^{old} ;

 Update the centered coordinates using (19) ;

M-scale-step: Evaluate s^{new} using (21);

M-rotation-step: Estimate \mathbf{R}^{new} via constrained non-linear optimization of (13) using (20) ;

M-covariance-step: Evaluate Σ^{new} using (23);

M-prior-step: Evaluate p^{new} using (22);

$\theta^{\text{old}} \leftarrow \theta^{\text{new}}$;

end

Optimal translation: Evaluate the translation vector using (9);

Result: Estimated scale s^* , rotation \mathbf{R}^* , translation \mathbf{t}^* , prior p^* , covariance Σ^* , and posterior probabilities of landmarks $\alpha_{1:N}$.

Algorithm 1: GUM-EM for robust estimation of the rigid transformation between two point sets.

4.3. Generalized Student Model

Another way to enforce robustness is to use the *generalized Student's t-distribution*, also known as the Pearson type VII distribution (Sun et al, 2010):

$$P(\mathbf{r}; \Sigma, \mu, \nu) = \int_0^\infty \mathcal{N}(\mathbf{r}; 0, w^{-1} \Sigma) \mathcal{G}(w, \mu, \nu) dw \\ = \frac{\Gamma(\mu + \frac{3}{2})}{|\Sigma|^{\frac{1}{2}} \Gamma(\mu) (2\pi\nu)^{\frac{3}{2}}} \left(1 + \frac{\|\mathbf{r}\|_{\Sigma}^2}{2\nu} \right)^{-(\mu + \frac{3}{2})} \quad (24)$$

where μ and ν are the parameters of the *prior gamma distribution* of w and $\Gamma(\cdot)$ is the gamma function. The distribution (24) differs from the standard Student's t-distribution in that the weight variable W , or the precision, is drawn from a gamma distribution with parameters μ and ν , instead of $\nu/2$ and $\nu/2$. Notice that in (24) ν and Σ appear only through their product, which means that an additional constraint is required to make the parameterization unique. One possibility is to constrain the determinant of the covariance, e.g. $|\Sigma| = 1$, which is equivalent to have an unconstrained Σ with $\nu = 1$. Unconstrained parameters are easier to deal with in inference algorithms. Therefore, we will rather assume without loss of generality that $\nu = 1$.

Notice that the posterior distribution of w is also a gamma distribution, namely the *posterior gamma distribution*:

$$P(w_n|\mathbf{r}_n; \Sigma, \mu, \nu) = \mathcal{N}(\mathbf{r}_n; 0, w_n^{-1} \Sigma) \mathcal{G}(w_n, \mu, \nu) \\ = \mathcal{G}(w_n; a, b_n), \quad (25)$$

with parameters:

$$a = \mu + \frac{3}{2}, \quad b_n = 1 + \frac{1}{2} \|\mathbf{r}_n\|_{\Sigma}^2. \quad (26)$$

²Note that the translation vector \mathbf{t} is evaluated outside the EM procedure.

Data: Centered point coordinates, i.e. (10). ;

Initialization of θ^{old} : Use the closed-form solution (Horn, 1987) to evaluate s^{old} and \mathbf{R}^{old} ; evaluate Σ^{old} using (15). Provide μ^{old} ;

while $\|\theta^{\text{new}} - \theta^{\text{old}}\| > \epsilon$ **do**

E-step: evaluate a^{new} and $b_{1:N}^{\text{new}}$ using (26) with θ^{old} , then evaluate $\bar{w}_{1:N}^{\text{new}}$ using (27) ;

Update the centered coordinates using (19), where $\alpha_{1:N}$ are replaced with $\bar{w}_{1:N}$;

M-scale-step: Evaluate s^{new} using (21);

M-rotation-step: Estimate \mathbf{R}^{new} with (13), (20) ;

M-covariance-step: Evaluate Σ^{new} using (29) ;

M-mu-step: Evaluate μ^{new} using (30) ;

$\theta^{\text{old}} \leftarrow \theta^{\text{new}}$;

end

Optimal translation: Evaluate the translation vector using (9);

Result: Estimated scale s^* , rotation \mathbf{R}^* , translation \mathbf{t}^* , covariance Σ^* , and landmark weights $w_{1:N}$.

Algorithm 2: The GStudent-EM for robust estimation of the rigid transformation between two point sets.

The posterior mean of the weight variable is:

$$\bar{w}_n = E[w_n | \mathbf{r}_n] = \frac{a}{b_n}. \quad (27)$$

As with the Gaussian-uniform model, we need to minimize the expected complete-data negative log-likelihood, $E_W[-\log P(\mathbf{r}_{1:N}, W_{1:N} | \mathbf{r}_{1:N}; \theta)]$ and in this case the parameter vector is $\theta = \{s, \mathbf{R}, \Sigma, \mu\}$ since we set $\nu = 1$. This yields the minimization of:

$$Q(\theta | (\mathbf{x}', \mathbf{y}')_{1:N}) = \frac{1}{2} \sum_{n=1}^N (\bar{w}_n \|\mathbf{y}'_n - s\mathbf{R}\mathbf{x}'_n\|_{\Sigma}^2 + \log |\Sigma|), \quad (28)$$

thus replacing $\alpha_{1:N}$ with $w_{1:N}$ in (19) and (20) to estimate the optimal rotation (13) and scale (21). The covariance matrix is estimated with:

$$\Sigma = \frac{1}{N} \sum_{n=1}^N \bar{w}_n (\mathbf{y}'_n - s\mathbf{R}\mathbf{x}'_n)(\mathbf{y}'_n - s\mathbf{R}\mathbf{x}'_n)^{\top} \quad (29)$$

The parameter μ is updated by solving the following equation, where $\Psi(\cdot)$ is the digamma function:

$$\mu = \Psi^{-1} \left(\Psi(a) - \frac{1}{n} \sum_{n=1}^N \log b_n \right). \quad (30)$$

We refer to this model as the *generalized Student* (GStudent) and the associated EM algorithm is summarized in Algorithm 2.

4.4. Algorithm Implementation and Analysis

Algorithm 1 and Algorithm 2 are expectation maximization (EM) procedures and it is well known that they have good convergence properties. One should notice that all the computations inside these algorithms are in closed-form, with the notable exception of the estimation of the rotation matrix. The

latter is parameterized with a unit quaternion which may be estimated via nonlinear constrained optimization. The unit-quaternion parameterization of rotations, i.e. Appendix Appendix A, has several advantages: (i) the number of parameters to be estimated is reduced from nine to four, (ii) the number of nonlinear constraints is reduced from seven constraints (six quadratic constraints, i.e. $\mathbf{R}^{\top} \mathbf{R} = \mathbf{I}$, and one quartic constraint, i.e. $|\mathbf{R}| = 1$) to one quadratic constraint ($\mathbf{q}^{\top} \mathbf{q} = 1$), (iii) the initialization is performed with the closed-form solution of (Horn, 1987) that uses a unit quaternion as well.

In practice, the constrained nonlinear optimization problem (13) is solved using the sequential quadratic programming method (Bonnans et al, 2006), more precisely a sequential least squares programming (SLSQP) solver³ is used in combination with a root-finding software package (Kraft, 1988). The SLSQP minimizer found at the previous EM iteration is used as an initial estimate at the current EM iteration. The closed-form method of (Horn, 1987) (please consult Appendix Appendix A) is used to initialize the unit-quaternion parameters at the start of the EM algorithm.

Other closed-form methods commonly used in computer vision, e.g. (Horn et al, 1988; Arun et al, 1987), perform singular value decomposition to extract an orthogonal matrix from the measurement matrix, but without the guarantee that the estimated matrix is a rotation, i.e. its determinant must be equal to +1. Appendix Appendix A summarizes the unit-quaternion closed-form method, which is based on estimating the smallest eigenvalue and eigenvector pair of a 4×4 semi-definite positive symmetric matrix – a well known mathematical problem yielding a straightforward numerical solver.

5. Analyzing the Robustness of Rigid Mapping

In order to quantify the performance of the proposed robust rigid-mapping algorithms, we devised an experimental protocol on the following grounds. Let $\mathbf{x}_{1:N}$ be a set of landmarks associated with the frontal view of a face. The set $\mathbf{y}_{1:N}^m$ is generated with:

$$\mathbf{y}_n^m(b) = s^m \mathbf{R}^m \mathbf{x}_n + \mathbf{t}^m + \mathbf{r}_n^m(b), \forall n \in \{1, \dots, N\}, \quad (31)$$

where $b > 0$ is a scalar that controls the level of noise and m is the trial index. As described in detail below, the noise level, b can be the variance of Gaussian isotropic noise, the total variance of Gaussian anisotropic noise, or the volume of uniformly distributed noise. The landmark coordinates are normalized such that $\forall n, \mathbf{x}_n \in [0, 1]^3$. For each noise level, we randomly generate M trials, namely M rigid mappings and M sets of N residuals $\mathbf{r}_{1:N} = \{\mathbf{r}_n\}_{n=1}^{n=N}$. For each trial m we estimate the rigid mapping parameters, $s^m, \mathbf{R}^m, \mathbf{t}^m$, and we measure the *root mean square error* (RMSE) between these estimated parameters and

³<https://docs.scipy.org/doc/scipy/reference/optimize.html>

the ground-truth parameters, $\tilde{s}^m, \tilde{\mathbf{R}}^m, \tilde{\mathbf{t}}^m$, namely:

$$E_s = \left(1/M \sum_{m=1}^M |s^m - \tilde{s}^m|^2\right)^{1/2}, \quad (32)$$

$$E_t = \left(1/M \sum_{m=1}^M \|\mathbf{t}^m - \tilde{\mathbf{t}}^m\|^2\right)^{1/2}, \quad (33)$$

$$E_{\mathbf{R}} = \left(1/M \sum_{m=1}^M \|\mathbf{R}^m - \tilde{\mathbf{R}}^m\|^2\right)^{1/2}, \quad (34)$$

The ground-truth rigid-mapping parameters are generated in the following way. For each trial m , the scale and the translation vector are generated from uniform distributions, namely $s^m \sim \mathcal{U}(0.5, 2)$ and $\mathbf{t}^m \sim \mathcal{U}(0.5, 5)^3$. The rotation matrix is parameterized by the pan, tilt and yaw angles, namely:

$$\mathbf{R} = \mathbf{R}_\gamma \mathbf{R}_\phi \mathbf{R}_\psi =$$

$$\begin{pmatrix} \cos \gamma & -\sin \gamma & 0 \\ \sin \gamma & \cos \gamma & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \cos \phi & 0 & \sin \phi \\ 0 & 1 & 0 \\ -\sin \phi & 0 & \cos \phi \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & \cos \psi & -\sin \psi \\ 0 & \sin \psi & \cos \psi \end{pmatrix}$$

A rotation matrix is obtained by randomly generating the pan, tilt and yaw angles, γ^m, ϕ^m, ψ^m , from a uniform distribution, $\mathcal{U}(-90^\circ, +90^\circ)$.

In order to generate residuals, $\mathbf{r}_{1:N}$, we simulate three types of noise:

- *Isotropic Gaussian noise:* $\mathbf{r} \sim \mathcal{N}(\mathbf{0}, \sigma \mathbf{I})$;
- *Anisotropic Gaussian noise:* $\mathbf{r} \sim \mathcal{N}(\mathbf{0}, \Sigma)$, and
- *Uniform noise:* $\mathbf{r} \sim \mathcal{U}(-a/2, a/2)^3$.

In the case of anisotropic noise, a covariance matrix must be randomly generated for each trial. This is done in the following way. Let $\Sigma = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$, with $\mathbf{Q} \in O(3)$ (an orthogonal matrix) and with $\mathbf{\Lambda} = \text{Diag}(\lambda_1, \lambda_2, \lambda_3)$, where the eigenvalues correspond to the variances along the eigenvectors – the directions of maximum variance. Let $\lambda = \lambda_1 + \lambda_2 + \lambda_3$ denote the total variance. A sample covariance matrix Σ is simulated by randomly generating an orthogonal matrix \mathbf{Q} and by randomly generating the three eigenvalues from a uniform distribution, $\mathcal{U}(0, 1)$.

We tested the following rigid mapping models and associated algorithms:

- *Horn:* Gaussian distribution with isotropic covariance, (Horn, 1987) and Appendix A;
- *Gen-Horn:* Gaussian distribution with anisotropic covariance, Section 4.1;
- *GUM-EM:* Gaussian-uniform mixture distribution, Algorithm 1, and
- *GStudent-EM:* Generalized Student's t-distribution, Algorithm 2.

The experiments were conducted in the following way. For each noise level, we simulated $M = 500$ trials for which we computed the RMSEs, namely eqs. (32), (33), and (34). For

each trial m we split the landmarks into an inlier set and an outlier set and the $N = 68$ landmarks are randomly assigned to one of these sets. The first experiment determines the percentage of outliers that can be handled by the robust algorithms, Figure 2. For this purpose, the percentage of outliers is increased from 10% to 60%. The inlier noise is drawn from an anisotropic Gaussian distribution with a total variance $\lambda = 0.0025$. The outlier noise is drawn from a uniform distribution with amplitude $a = 1.5$ (remember that the landmark coordinates are normalized to lie in the interval $[0, 1]$). The curves plotted in Figure 2 show that the RMSE associated with non robust methods, i.e. Horn and Gen-Horn increase monotonically. On the contrary, the robust algorithms, GUM-EM and GStudent-EM, have a radically different behavior. After a short increase, the RMSE remains constant, and then it increases again.

In the other experiments, the number of inliers was set to be equal to the number of outliers and we experimented with the three noise types already mentioned. Figure 3 shows the RMSEs when inlier noise is drawn from an isotropic Gaussian distribution with $\sigma = 0.0025$, while outlier noise is drawn from a uniform distribution whose volume is increased from $a = 0.2$ to $a = 1.0$. Similarly, Figure 5 shows the RMSEs for the case when inlier noise is drawn from an anisotropic Gaussian distribution with total variance $\lambda = 0.0025$, while outlier noise is drawn from a uniform distribution whose volume is increased from $a = 0.2$ to $a = 1.0$. Finally, Figure 5 shows the RMSEs when inlier noise is drawn from an anisotropic Gaussian distribution with total variance $\lambda = 0.0025$, while outlier noise is drawn from an anisotropic Gaussian distribution with total variance varying from $\lambda = 0.2$ to $\lambda = 1.0$.

These experiments clearly show that the two classes of methods (non-robust and robust) behave differently. The performance of non-robust rigid mapping decreases monotonically in the presence of outliers with increasing levels of noise. The robust methods can deal with up to 50% of outliers affected by a substantial noise level (1.5 times the size of the image). There is no evidence that the Gen-Horn algorithm performs better than the standard Horn algorithm. Nevertheless, Gen-Horn provides interesting information about the 3D structure of the estimated anisotropic covariance. The GUM-EM algorithm performs slightly better than the GStudent-EM algorithm, in particular in the presence of outliers drawn from a uniform distribution.

6. Measuring the Performance of 3D Face Alignment

In this section we describe an unsupervised methodology for quantitatively assessing the performance of 3DFA algorithms. The idea of the proposed benchmarking is to apply 3DFA to a dataset of face images in order to extract 3D landmarks, to robustly estimate the rigid transformation that maps these facial landmarks into a 3D landmark model, and to analyze the discrepancy between the extracted 3D landmarks and the model. Based on a confidence score, it is then possible to decide whether a landmark is correctly localized or not. This allows to assess the overall performance of a 3DFA algorithm as

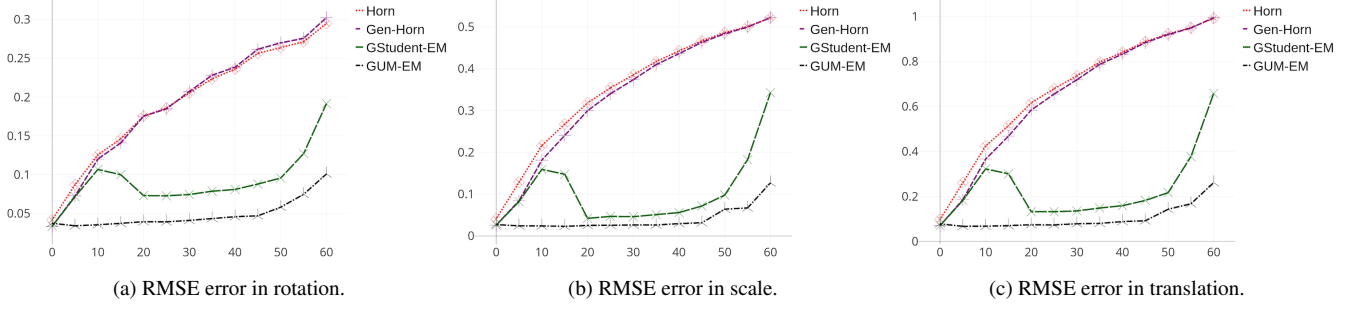


Fig. 2: RMSE error as a function of the percentage of outliers: inliers are affected by anisotropic Gaussian noise with total variance $\lambda = 0.0025$, while the percentage of outliers, affected by uniform noise with amplitude $a = 1.5$, increases from 0% to 60%.

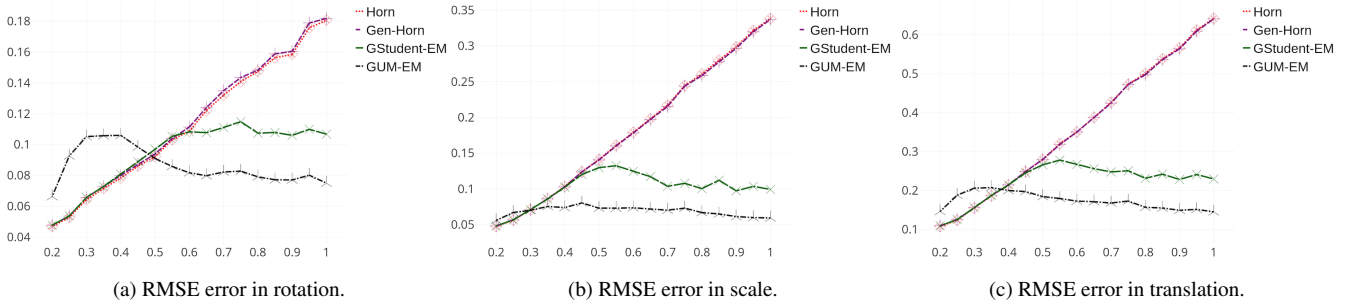


Fig. 3: RMSE error as a function of uniform noise affecting a fixed number of outliers: inliers (50%) are affected by isotropic Gaussian noise with variance $\sigma = 0.0025$, while outliers (50%) are affected by uniform noise of increasing amplitude $a \in [0.2, 1.0]$.

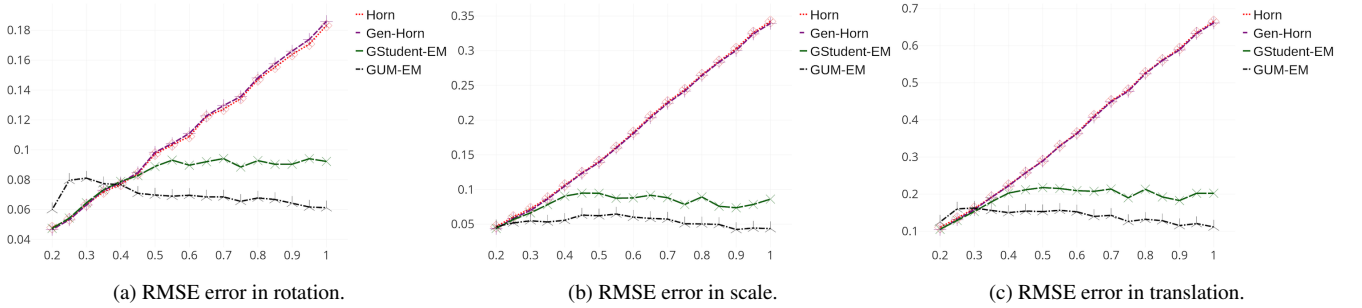


Fig. 4: RMSE error as a function of uniform noise affecting a fixed number of outliers: inliers (50%) are affected by anisotropic Gaussian noise with total variance $\lambda = 0.0025$, while outliers (50%) are affected by uniform noise of increasing amplitude $a \in [0.2, 1.0]$.

well as its behavior with respect to various perturbations, such as occlusions or motion blur.

6.1. Neutral Frontal Landmark Model

We start by computing a *neutral frontal landmark model* $\mathbf{y}_{1:N}$ in the following way. For this purpose, we use a dataset \mathcal{D}_1 of K images of neutral faces (frontal viewing, no expression and no interfering object causing occlusion) and we extract N landmarks from each one of these K faces, $\{\mathbf{y}_{1:N,k}\}_{k=1}^K$. Then we use the landmark coordinates to compute the directions of maximum variance (or the principal components) of each face. By

aligning these directions over the dataset, we compute a mean for each landmark, namely

$$\mathbf{y}_n = 1/K \sum_{k=1}^K \mathbf{y}_{n,k}, \forall n, 1 \leq n \leq N. \quad (35)$$

6.2. Statistical Frontal Landmark Model

We now explain how a *statistical frontal landmark model* is built, namely $\{\mathbf{p}_{1:N}, \mathbf{C}_{1:N}\}$, where $\mathbf{p}_{1:N}$ is the set of N means and $\mathbf{C}_{1:N}$ is the set of N covariance matrices associated with the statistical frontal landmark model. For this purpose, we use

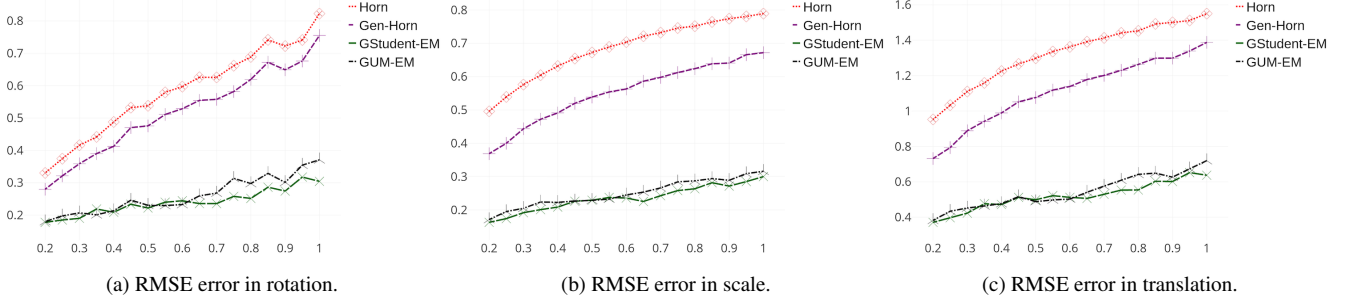


Fig. 5: RMSE error as a function of anisotropic Gaussian noise affecting a fixed number of outliers: inliers (50%) are affected by anisotropic Gaussian noise with total variance $\lambda = 0.0025$, while outliers (50%) are affected by anisotropic Gaussian noise with total variance $\lambda \in [0.2, 1.0]$.

another dataset \mathcal{D}_2 that contains L images of faces with the following characteristics: arbitrary poses, arbitrary expressions, both speaking and silent faces, but with no external sources of perturbation such as the presence of interfering object that may cause occlusions. We extract 3D landmarks from these L images using a 3DFA algorithm, namely $\{x_{1:N,l}\}_{l=1}^L$, and we use either GUM-EM (Algorithm 1) or GStudent-EM (Algorithm 2) to robustly estimate the rigid transformations between each landmark-set l and the neutral frontal landmark-set $x_{1:N,l}$ and $y_{1:N}$. Based on this, we obtain L rigid-mapping parameters (one for each l): L scale factors, L rotations and L translations: $\{s_l^{\text{Alg}}, \mathbf{R}_l^{\text{Alg}}, \mathbf{t}_l^{\text{Alg}}\}_{l=1}^L$, where the over-script Alg denotes a robust algorithm, namely either GUM-EM or GStudent-EM. We remind that both algorithms provide a figure of merit for each landmark: posterior probabilities $\{\alpha_{n,l}\}_{n=1}^{n=N}$ in the case of GUM-EM, i.e. (18), and precision means $\{\bar{w}_{n,l}\}_{n=1}^{n=N}$ in the case of GStudent-EM, i.e. (27). Applying one of these robust rigid-alignment methods provides frontal landmarks, $\{\tilde{x}_{1:N,l}^{\text{Alg}}\}_{l=1}^L$, namely:

$$\tilde{x}_{n,l}^{\text{Alg}} = s_l^{\text{Alg}} \mathbf{R}_l^{\text{Alg}} x_{n,l} + \mathbf{t}_l^{\text{Alg}}. \quad (36)$$

There are two different expressions for the posterior means and posterior covariances for GUM-EM and for GStudent-EM, respectively:

$$\mathbf{p}_n^{\text{GUM}} = \frac{\sum_{l=1}^L \alpha_{n,l} \tilde{x}_{n,l}^{\text{GUM}}}{\sum_{l=1}^L \alpha_{n,l}}, \quad (37)$$

$$\mathbf{C}_n^{\text{GUM}} = \frac{\sum_{l=1}^L \alpha_{n,l} (\tilde{x}_{n,l}^{\text{GUM}} - \mathbf{p}_n^{\text{GUM}})(\tilde{x}_{n,l}^{\text{GUM}} - \mathbf{p}_n^{\text{GUM}})^\top}{\sum_{l=1}^L \alpha_{n,l}}, \quad (38)$$

and

$$\mathbf{p}_n^{\text{GSt}} = \frac{\sum_{l=1}^L \bar{w}_{n,l} \tilde{x}_{n,l}^{\text{GSt}}}{\sum_{l=1}^L \bar{w}_{n,l}}, \quad (39)$$

$$\mathbf{C}_n^{\text{GSt}} = \frac{\sum_{l=1}^L \bar{w}_{n,l} (\tilde{x}_{n,l}^{\text{GSt}} - \mathbf{p}_n^{\text{GSt}})(\tilde{x}_{n,l}^{\text{GSt}} - \mathbf{p}_n^{\text{GSt}})^\top}{\sum_{l=1}^L \bar{w}_{n,l}}. \quad (40)$$

Notice that (39) and (40) compute a mean and a covariance for landmark n over the entire dataset. Hence, and unlike in (29), the covariance should be normalized with the sum of the weights.

6.3. Unsupervised Confidence Test

We now develop an unsupervised (statistical) confidence test for assessing whether the accuracy of a landmark, i.e. its 3D coordinates, is within (inlier) or outside (outlier) an expected range (Savage, 1972). Let us drop the algorithm over-script and let $\mathbf{C}_n = \mathbf{Q}_n \mathbf{\Lambda}_n \mathbf{Q}_n^\top$ be the eigen factorization of \mathbf{C}_n , where \mathbf{Q}_n is an orthonormal matrix and $\mathbf{\Lambda}_n$ is a diagonal matrix containing the eigenvalues. We can now project each landmark (n, l) onto the space spanned by the three eigenvectors of this matrix:

$$\tilde{z}_{n,l} = \mathbf{Q}_n^\top (\tilde{x}_{n,l} - \mathbf{p}_n) \quad (41)$$

Landmark (n, l) is an inlier with 99% confidence if $\tilde{z}_{n,l}$ lies inside the ellipsoid whose half-axes are three times the standard deviations, or $3\sqrt{\lambda_n^1}$, $3\sqrt{\lambda_n^2}$, $3\sqrt{\lambda_n^3}$, where $\{\lambda_n^1, \lambda_n^2, \lambda_n^3\}$ are the eigenvalues of \mathbf{C}_n , or

$$\tilde{z}_{n,l}^\top \tilde{\mathbf{\Lambda}}_n^{-1} \tilde{z}_{n,l} \leq 1 \quad (42)$$

where $\tilde{\mathbf{\Lambda}}_n = 9\mathbf{\Lambda}_n$. Combining (41) and (42), yields $(\tilde{x}_{n,l} - \mathbf{p}_n)^\top \mathbf{Q}_n \tilde{\mathbf{\Lambda}}_n^{-1} \mathbf{Q}_n^\top (\tilde{x}_{n,l} - \mathbf{p}_n) \leq 1$. With the notation

$$\tilde{\mathbf{C}}_n = \mathbf{Q}_n \tilde{\mathbf{\Lambda}}_n \mathbf{Q}_n^\top. \quad (43)$$

The 99% confidence test writes:

$$\begin{cases} \text{if } \|\tilde{x}_{n,m} - \mathbf{p}_n\|_{\tilde{\mathbf{C}}_n} \leq 1 & (n, m) = \text{inlier} \\ \text{otherwise} & (n, m) = \text{outlier} \end{cases} \quad (44)$$

Based on this confidence test, we can now build a *confidence-test accuracy* (the higher the better) associated with a sample face m , namely:

$$u(m) = \frac{1}{N} \sum_{n=1}^N \mathcal{I}(\|\tilde{x}_{n,l} - \mathbf{p}_n\|_{\tilde{\mathbf{C}}_n} \leq 1), \quad (45)$$

where $\mathcal{I}(\cdot)$ denotes the indicator function. Notice that (45) corresponds to the percentage of inliers, i.e. landmarks that, once scaled, rotated and translated, lie inside the confidence volume. Therefore, (45) can be used to assess whether the pose has been correctly estimated, namely $u \leq 50\%$, or not (please consult Section 5). For a test dataset \mathcal{D}_3 composed of M samples, one can then compute the *mean confidence test accuracy* (the higher the better):

$$\mathbf{U} = \frac{1}{M} \sum_{m=1}^M u(m) \quad (46)$$

6.4. Correlation with Supervised Metrics

In general, datasets of faces come with their ground-truth annotations, and we denote with $\hat{\mathbf{x}}_{1:N,1:M}$ the set of ground-truth landmarks associated with the dataset \mathcal{D}_2 . We modify (1) to be able to build a metric that counts the proportion of inliers, namely the *ground-truth accuracy* (the higher the better):

$$s(m) = \frac{1}{N} \sum_{n=1}^N \mathcal{I}(\|\mathbf{x}_{n,m} - \hat{\mathbf{x}}_{n,m}\|/d_m \leq \varepsilon), \quad (47)$$

where ε is a user-defined threshold that corresponds to the quality of the ground-truth landmarks. Based on this we can compute the *mean ground-truth accuracy* (S):

$$S = \frac{1}{M} \sum_{m=1}^M s(m) \quad (48)$$

Finally, another interesting metric is the correlation coefficient between the above unsupervised and supervised metrics:

$$\text{Cor} = \frac{\sum_{m=1}^M (u(m) - U)(s(m) - S)}{\left(\sum_{m=1}^M (u(m) - U)^2 \sum_{m=1}^M (s(m) - S)^2\right)^{1/2}} \quad (49)$$

7. Experimental Results

7.1. Neutral Frontal Landmark Model

The neutral frontal landmark model was trained in-the-wild by harvesting web images and using a face detector and a head-pose estimator in order to select frontal faces. These images were visually inspected to guarantee shape and aspect variabilities as well as neutral facial expressions. This process yields a dataset \mathcal{D}_1 composed of 1,000 images. We used the 3DFA method of (Feng et al, 2018) to extract landmarks from each face in the dataset. Next, we aligned them (please consult Section 6.1) and computed the landmark means using (35). Figure 6 shows a few examples of images from this dataset as well as the detected landmarks. Figure 7 show the neutral frontal landmark model thus obtained.

7.2. Statistical Frontal Landmark Model

The statistical frontal landmark model was trained from the YawDD dataset (Abtahi et al, 2014). This dataset contains 322 videos which is equivalent to approximatively 300,000 images. The face images in this dataset have large variabilities in terms of face shapes, face aspects, head poses and facial expressions. All the images were processed with no human intervention, namely: face detection, 3D face alignment, and robust rigid alignment with the neutral face landmarks just described. This yields the statistical face landmark model described in Section 6.2. For that purpose we used two 3DFA methods and the two robust alignment algorithms described in this paper. Hence, there are four possible 3DFA and robust alignment combinations that we used to train four different models:

- 3DFA1/GUM-EM: (Bulat and Tzimiropoulos, 2016) and GUM-EM (Algorithm 1),
- 3DFA2/GUM-EM: (Feng et al, 2018) and GUM-EM (Algorithm 1),
- 3DFA1/GStudent-EM: (Bulat and Tzimiropoulos, 2016) and GStudent-EM (Algorithm 2), and
- 3DFA2/GStudent-EM: (Feng et al, 2018) and GStudent-EM (Algorithm 2).

Figure 8 shows the statistical frontal landmark models obtained with these four combination. In this figure, the dots correspond to the posterior means, i.e. (37) and (39), while the ellipses correspond to image projections of the ellipsoids defined by (43).

7.3. Performance Evaluation of 3D Face Alignment

Once the neutral-frontal and statistical-frontal models are computed using datasets \mathcal{D}_1 and \mathcal{D}_2 , respectively, we use a third dataset, \mathcal{D}_3 , to empirically assess the performance of four 3DFA algorithms using the unsupervised confidence test introduced in Section 6.3:

- 3DFA1 (Bulat and Tzimiropoulos, 2016) (ECCVW'16) and (Bulat and Tzimiropoulos, 2017) (ICCV'17);
- 3DFA2 (Feng et al, 2018) (ECCV'18);
- 3DFA3 (Zhu et al, 2016) (CVPR'18) and (Zhu et al, 2019) (PAMI'19), and
- 3DFA4 (Tu et al, 2020) (IEEE TMM'20).

For this purpose we used AFLW2000-3D (Zhu et al, 2016) as a test dataset, consisting of 2000 images with large-pose variations. More precisely, the yaw angles (vertical axis of rotation) in the following intervals $[0^\circ, \pm 30^\circ]$ for 1306 faces, in the interval $[\pm 30^\circ, \pm 60^\circ]$ for 462 faces and in the interval $[\pm 60^\circ, \pm 90^\circ]$ for 232 faces. The dataset contains a large variety of facial shapes, facial expressions and illuminations conditions. Moreover, there are many faces with partial occlusions caused by the presence of hair, hands, handheld objects, glasses, etc. Notice that large poses induce partial occlusions as well.

Each image in the AFLW2000-3D dataset is annotated with 68 3D landmarks. This semi-automatic annotation is performed by fitting a 3D deformable model to a dataset of 2D face images, e.g. (Ghiasi and Fowlkes, 2014). Nevertheless, as noted in (Bulat and Tzimiropoulos, 2017), many of the annotated landmarks in this dataset have large localization errors, especially in the case of profile views. Hence, performance evaluation based on supervised metrics are prone to errors. In (Bulat and Tzimiropoulos, 2017) it is visually shown that in these extreme poses their 3DFA method yields more precise landmark localization than the automatically annotated ones. Based on these observations, we applied our unsupervised performance analysis to the annotated landmarks as well, yielding the following combinations:

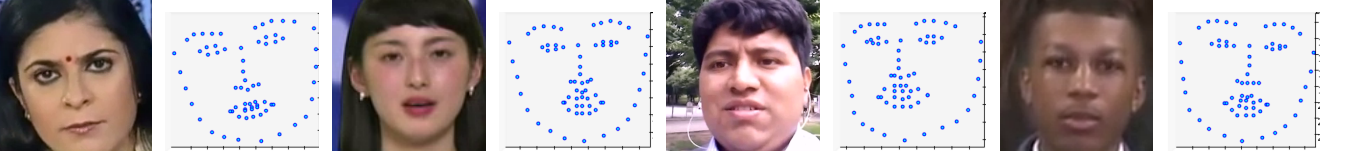


Fig. 6: Examples of faces and corresponding landmarks used to compute a neutral frontal landmark model.

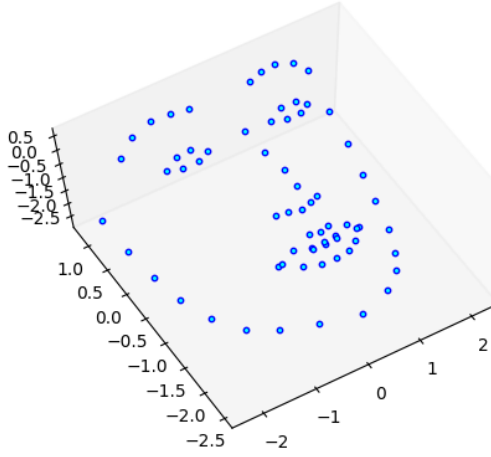


Fig. 7: A 3D view of the neutral frontal landmark model.

- GT/GUM-EM: Ground-truth landmarks provided by (Zhu et al, 2016) and GUM-EM, and
- GT/GStudent-EM: Ground-truth landmarks provided by (Zhu et al, 2016) and GStudent-EM.

The results based on computing the mean confidence-test accuracies, i.e. (46) are summarized in Table 1. We remind that we used different datasets for training the neutral and statistical face landmark models, i.e. \mathcal{D}_1 and \mathcal{D}_2 , and for assessing the performance of the various combinations of 3DFA methods and robust-rigid mappings, i.e. \mathcal{D}_3 . The means, evaluated over the confidence-test scores obtained with the annotated (ground-truth) landmarks (the last two rows of Table 1), are equal to 0.70 and to 0.65, respectively, which seems to confirm that the ground-truth landmark locations in the AFLW2000-3D dataset contain a substantial amount of errors and that, overall, both 3DFA methods that we analyzed, (Bulat and Tzimiropoulos, 2016) and (Feng et al, 2018), predict landmark locations that are more accurate than the ground-truth locations themselves.

We now compute correlation coefficients, i.e. (49), between the unsupervised and supervised metrics, i.e. (45), and (47). The results are reported in Table 2. Notice however that the supervised scores depend on the choice of the parameter ε . As done in (Bulat and Tzimiropoulos, 2017), this parameter was adjusted to eliminate samples yielding a very low score. With $\varepsilon = 0.1$ in (48) we obtained the following scores:

- 3DFA1: $S = 0.57$,

- 3DFA2: $S = 0.60$,
- 3DFA3: $S = 0.49$,
- 3DFA4: $S = 0.09$.

These scores are comparable with the scores reported in (Bulat and Tzimiropoulos, 2017) which uses a different normalization parameter. Notice that the scores obtained with the proposed confidence test, i.e. Table 1, are higher than these scores.

In the light of these results, we attempted to analyse the effect of eliminating inaccurate ground-truth landmark annotations from the benchmark just described. Let us define

$$\mathcal{M}_\tau = \{m \mid \hat{u}(m) \geq \tau\}, \quad \mathcal{M}_\tau \subset \mathcal{D}_3, \quad (50)$$

where $\hat{u}(m)$ denotes the value of the unsupervised score (45) associated with the ground-truth landmarks of face sample m . We see that when τ increases, the accuracy of the ground-truth landmark annotations contained in the subset \mathcal{M}_τ increases as well, at the price of drastically decreasing the number of inlying samples, which in turn lowers down the statistical significance of the resulting scores. The correlation coefficient is then computed with the following formula:

$$\text{Cor}(\tau) = \frac{\sum_{m \in \mathcal{M}_\tau} (u(m) - U)(s(m) - S)}{\left(\sum_{m \in \mathcal{M}_\tau} (u(m) - U)^2 \sum_{m \in \mathcal{M}_\tau} (s(m) - S)^2 \right)^{1/2}} \quad (51)$$

Figure 9-left shows the correlation as a function of τ where the 3DFA1/GStudent-EM method was used to build the confidence test, i.e. Figure 8-c, while Figure 9-right shows the corresponding p -value (the smaller the better) with a significance level of 0.005: a very small p -value is an indicator of the statistical significance of the correlation measure. The red dots in these plots correspond to a p -value not satisfying the significance level just mentioned.

In the light of these experiments, we conclude that the proposed methodology for assessing the performance of 3DFA methods is not biased by the quality of landmark annotation, whether the latter is automatic or human-assisted. The experiments suggest that the proposed unsupervised methodology could be used (i) to assess the quality of landmark annotation itself and (ii) to remove badly annotated landmarks.

We now illustrate the proposed performance analysis method with a few examples from the AFLW2000-3D. In all these examples, the statistical face model is trained with 3DFA1/GStudent-EM and GStudent-EM is used for testing. Figure 10 shows results obtained with 3DFA1 (Bulat and Tzimiropoulos, 2016), Figure 11 shows results obtained with

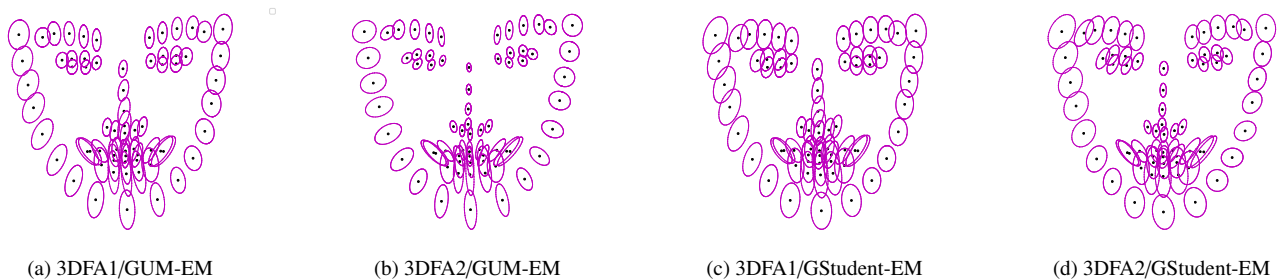


Fig. 8: Statistical frontal landmark models obtained with two 3DFA methods and with the proposed robust rigid-mapping algorithms.

3DFA2, Figure 12 shows results obtained with 3DFA3, and Figure 13 shows results obtained with 3DFA4.

We now show results obtained by applying GStudent-EM to the ground-truth landmarks associated with the AFLW2000-3D dataset (Zhu et al, 2016), or GT/GStudent-EM. Some best-score examples are shown in Figure 14 and some worse-score results are shown in Figure 15. Both the results reported in Table 1 and these examples suggest that the ground-truth annotations should be handled with caution.

We also gathered a dataset of 30 animal faces in order to analyse the performance of 3DFA with non-human faces. We processed these 30 images with the four 3DFA methods, as above, but only (Bulat and Tzimiropoulos, 2016) yielded exploitable results. As with human faces, for each one of these animal images we computed the percentage of inliers, i.e. (45). Figure 16 shows the six best and the six worst scores. The high scores observed on some of these animal faces correspond to failures of the proposed confidence test. In fact, the 3DFA method itself failed in these cases because it predicted a landmark pattern that corresponds to a human face. These results should, however, be interpreted with caution, because neither the 3DFA method of (Bulat and Tzimiropoulos, 2016) nor the proposed methodology were trained with animal faces.

8. Discussion and Conclusions

Landmarks predicted by face alignment methods, whether 2D or 3D, are inherently affected by noise and they contain outliers. We propose an unsupervised methodology to characterize the performance of 3DFA. First, we learn a statistical frontal landmark model, namely posterior mean and covariance for each 3D landmark. This pose-invariant model is materialized by an ellipsoidal volume of confidence that encapsulates variabilities due to face identity, face expression, occlusion, and detection noise. Second, the landmarks predicted by a 3DFA method are rigidly and robustly mapped onto this frontal model in order to be able to compare their locations with the model locations. A landmark that falls inside its corresponding ellipsoidal volume is labeled as inlier with 99% confidence, while a landmark that lies outside this volume is labeled as outlier. The ability (i) to separate pose from expression and from identity in a robust manner and (ii) to discriminate between inlying and outlying landmarks stands in contrast with existing evaluation

metrics that compute the mean distance between ground-truth landmark locations and predicted locations.

We note that none of the 3DFA evaluation metrics proposed so far exploit the concept of bringing all the facial landmarks into a canonical frontal pose. This is performed using a rigid alignment that is embedded into a maximum-likelihood estimator that uses a robust probability distribution function. This may well be viewed either as a robust pose estimator or as a mechanism that yields an expression- and identity-preserving frontal landmark representations. In turn, this enables temporal analysis of facial expressions and of lip movements.

When applied to the AFLW2000-3D dataset, the proposed analysis reveals that ground-truth landmark locations, provided by this annotated set of faces, contain 0.67 inliers, on an average (the last two rows of Table 1), which is less than the percentage of inliers associated with the three best-performing 3DFA methods that we analyzed. This result confirms the conclusions of (Bulat and Tzimiropoulos, 2017) that these annotations contain many large errors. To better understand these results, we also computed the correlation between the proposed unsupervised metric and the supervised metric. The correlation coefficients are in between 0.12 and 0.33 for the best-performing methods, as reported in Table 2. Interestingly, these correlation coefficients increase monotonically as the bad annotations are eliminated, which is illustrated in Figure 9. We conjecture that the proposed methodology can be used to eliminate annotation errors from a dataset of faces. Alternatively, it is also possible to automatically annotate a dataset using existing 3DFA methods and selecting the best localization for each landmark.

Because the proposed methodological pipeline makes use of 3DFA and of robust rigid mapping both for training and for testing, one may argue that the associated performance metric is biased. We empirically showed that the analysis is agnostic to various combinations of methods used for training the model and for the tests themselves. Moreover, we showed that the method could also be used to assess the quality of facial landmark annotations, in particular those annotations that are obtained automatically based on 3D deformable-model fitting. We therefore conclude that the interest of the proposed methodology is twofold, namely (i) to assess the accuracy, the repeatability and the reliability of the predictions obtained with 3DFA algorithms, and (ii) to evaluate the quality of the landmarks predicted from a test face.

Table 1: Performance analysis based on the proposed unsupervised metrics. The numbers correspond to the proportion of inliers (the higher the better) computed using (45) and (46) over a dataset that contains 2,000 face images and 68 landmarks per face.

Alignment method using dataset \mathcal{D}_3 :	Statistical face models trained with \mathcal{D}_1 and \mathcal{D}_2 datasets:				
	3DFA1/GUM-EM	3DFA2/GUM-EM	3DFA1/GStudent-EM	3DFA2/GStudent-EM	Mean
3DFA1/GUM-EM	0.89	0.65	0.93	0.80	0.82
3DFA2/GUM-EM	0.93	0.88	0.95	0.93	0.92
3DFA3/GUM-EM	0.98	0.96	0.98	0.98	0.98
3DFA4/GUM-EM	0.11	0.06	0.16	0.12	0.11
3DFA1/GStudent-EM	0.80	0.57	0.88	0.74	0.75
3DFA2/GStudent-EM	0.84	0.76	0.90	0.88	0.84
3DFA3/GStudent-EM	0.93	0.85	0.96	0.94	0.92
3DFA4/GStudent-EM	0.18	0.12	0.23	0.18	0.18
GT/GUM-EM	0.73	0.54	0.82	0.71	0.70
GT/GStudent-EM	0.67	0.48	0.78	0.66	0.65

Table 2: Correlation coefficients computed with (49) (the higher the better) between unsupervised and supervised metrics.

Alignment method using dataset \mathcal{D}_3 :	Statistical face models trained with \mathcal{D}_1 and \mathcal{D}_2 datasets:				
	3DFA1/GUM-EM	3DFA2/GUM-EM	3DFA1/GStudent-EM	3DFA2/GStudent-EM	Mean
3DFA1/GUM-EM	0.30	0.28	0.28	0.29	0.29
3DFA2/GUM-EM	0.33	0.37	0.30	0.33	0.33
3DFA3/GUM-EM	0.16	0.17	0.16	0.15	0.16
3DFA4/GUM-EM	0.05	0.04	0.05	0.04	0.04
3DFA1/GStudent-EM	0.25	0.26	0.26	0.26	0.26
3DFA2/GStudent-EM	0.25	0.28	0.23	0.26	0.25
3DFA3/GStudent-EM	0.13	0.13	0.13	0.12	0.13
3DFA4/GStudent-EM	0.05	0.04	0.07	0.03	0.13

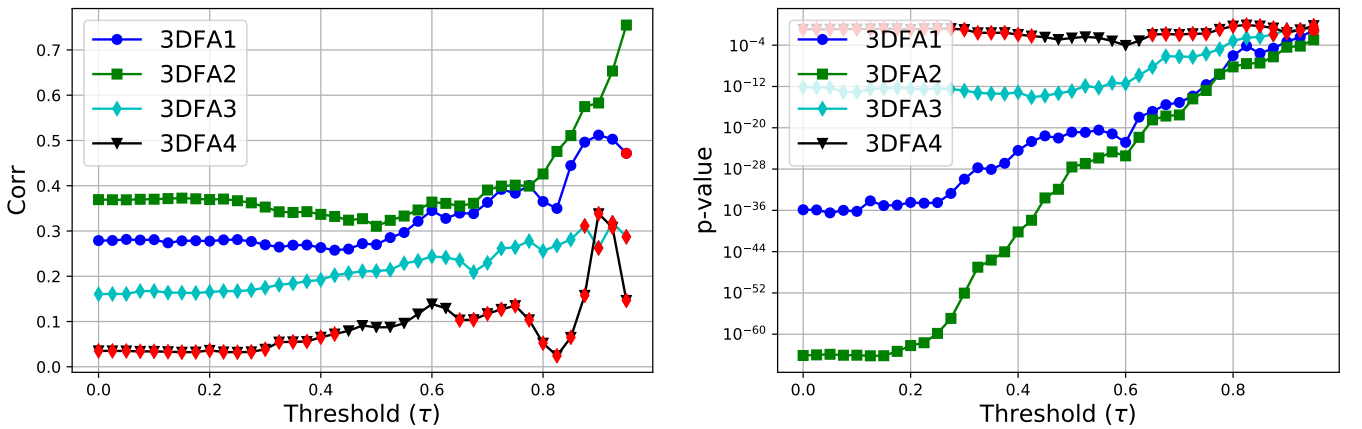


Fig. 9: Correlation between the unsupervised and supervised metrics (left) and corresponding p-value (right) as a function of the accuracy of the ground-truth landmark annotations. The red dots correspond to a p-value not satisfying a significance level, which is set to 0.005 in these experiments.

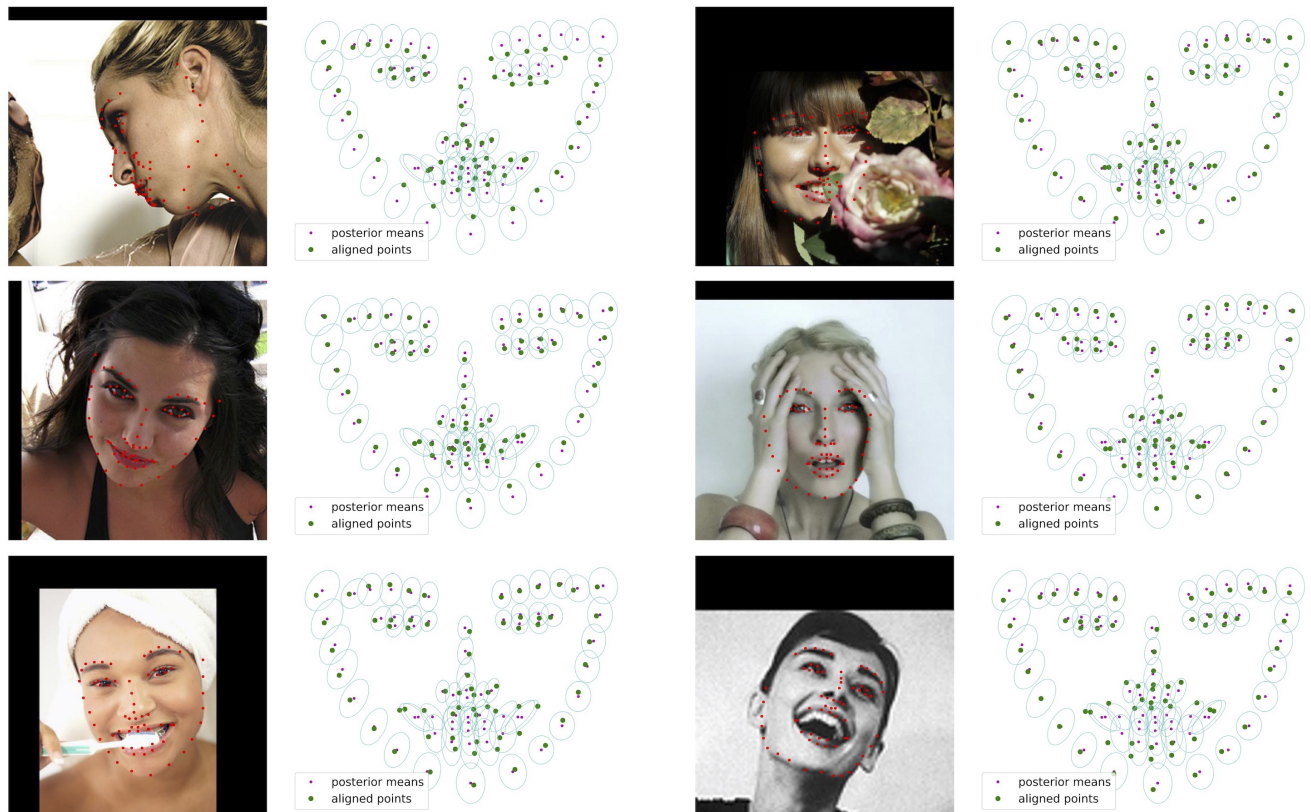


Fig. 10: A few examples obtained with 3DFA1 (Bulat and Tzimiropoulos, 2016) and with GStudent-EM.

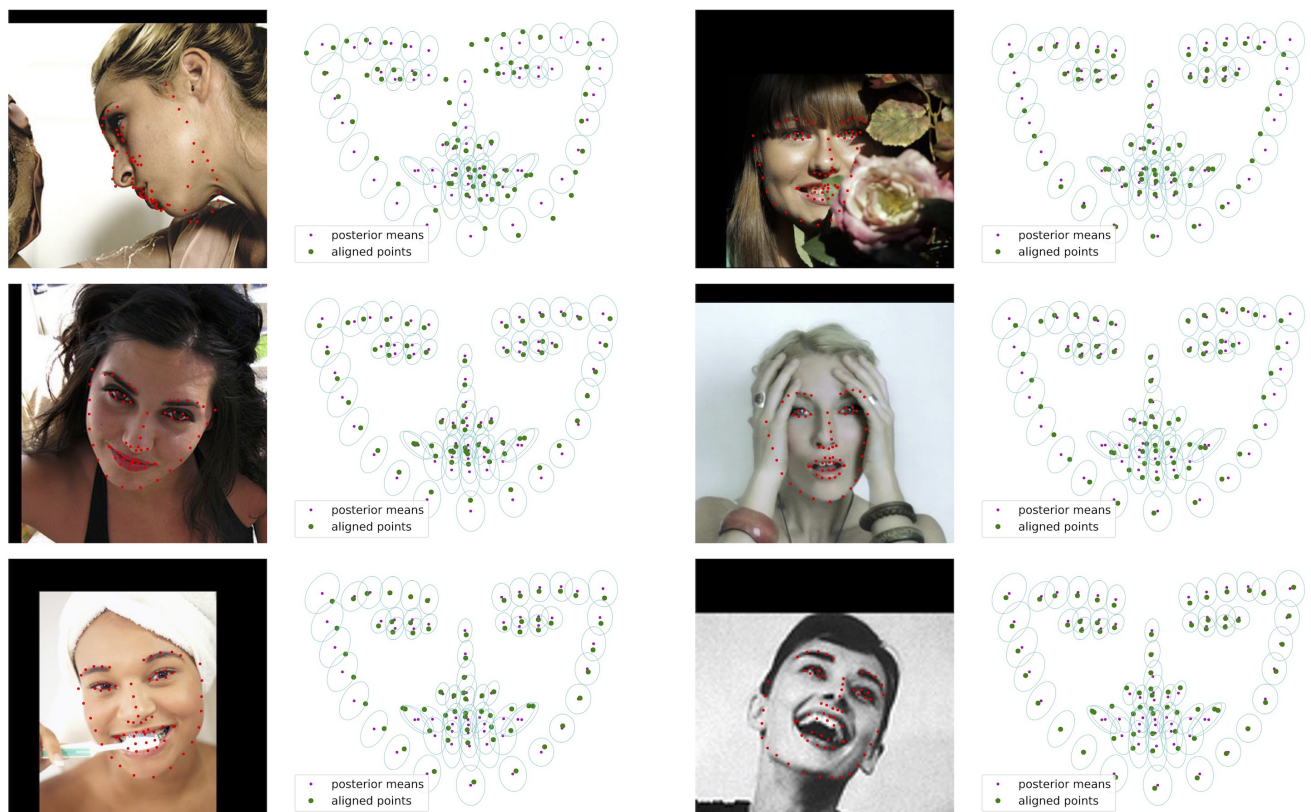


Fig. 11: A few examples obtained with 3DFA2 (Feng et al, 2018) and with GStudent-EM.

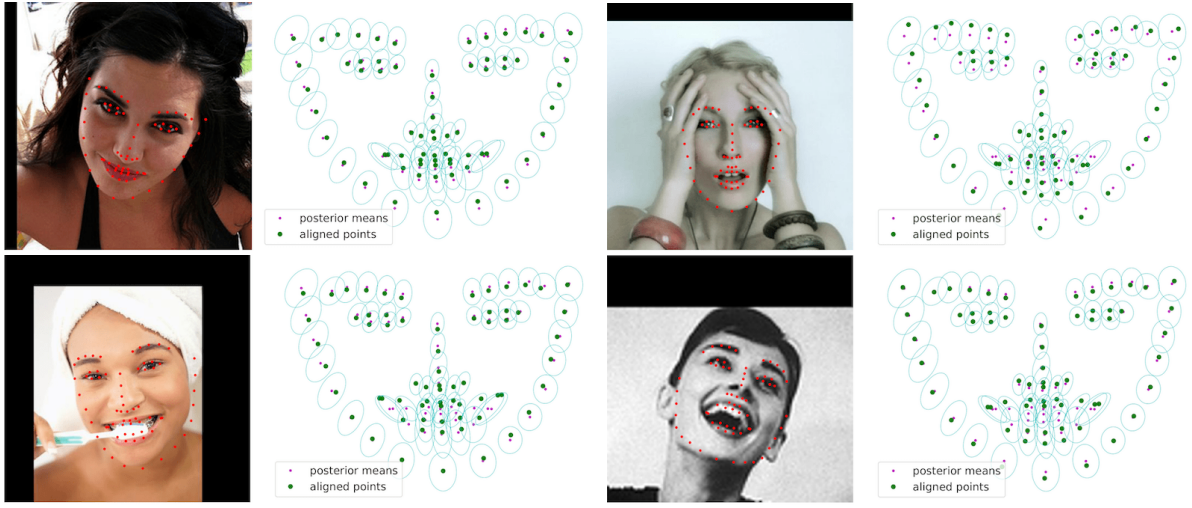


Fig. 12: A few examples obtained with 3DFA3 (Zhu et al, 2016) and with GStudent-EM.

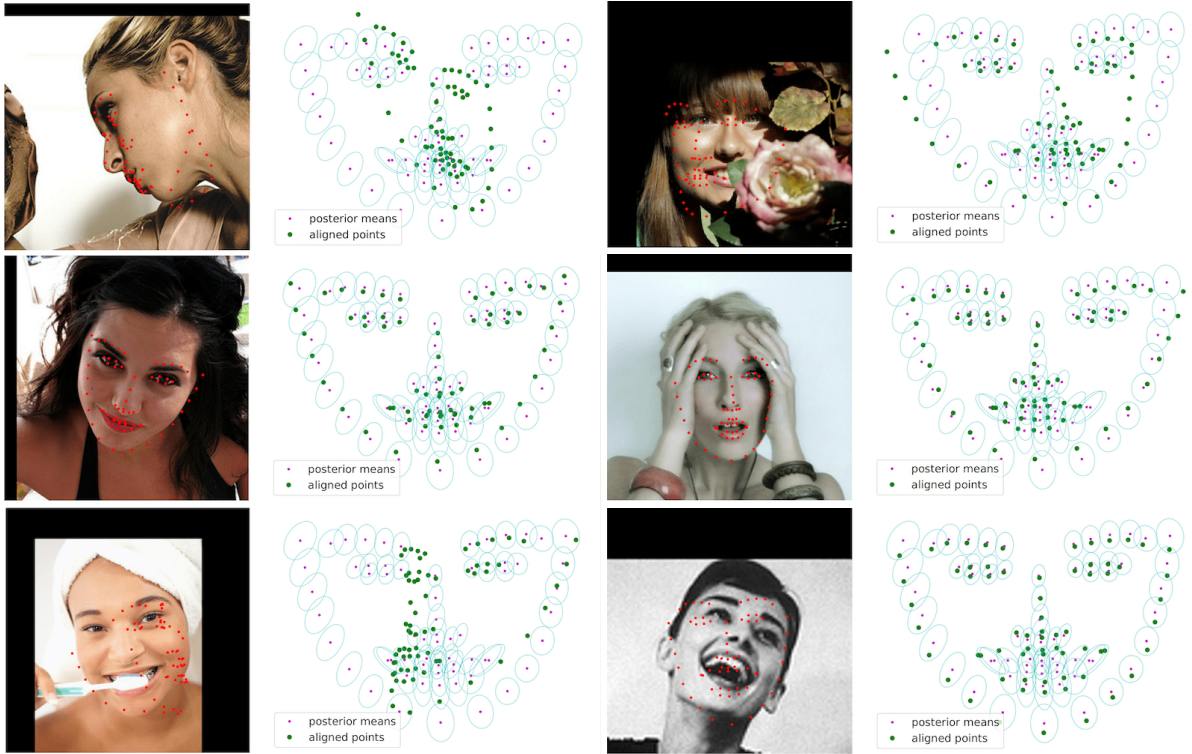


Fig. 13: A few examples obtained with 3DFA4 (Tu et al, 2020) and with GStudent-EM.

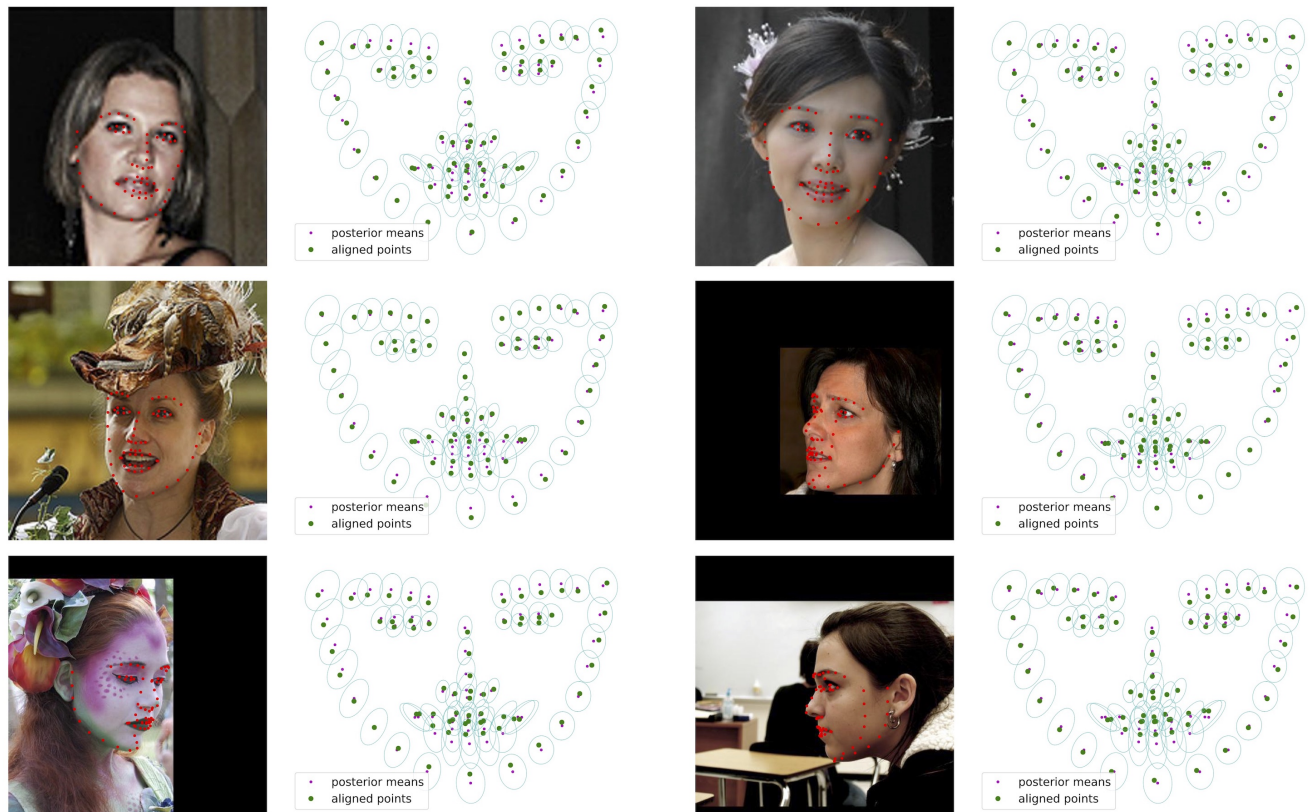


Fig. 14: Some examples of the best scores obtained with GStudent-EM and the ground-truth landmarks.

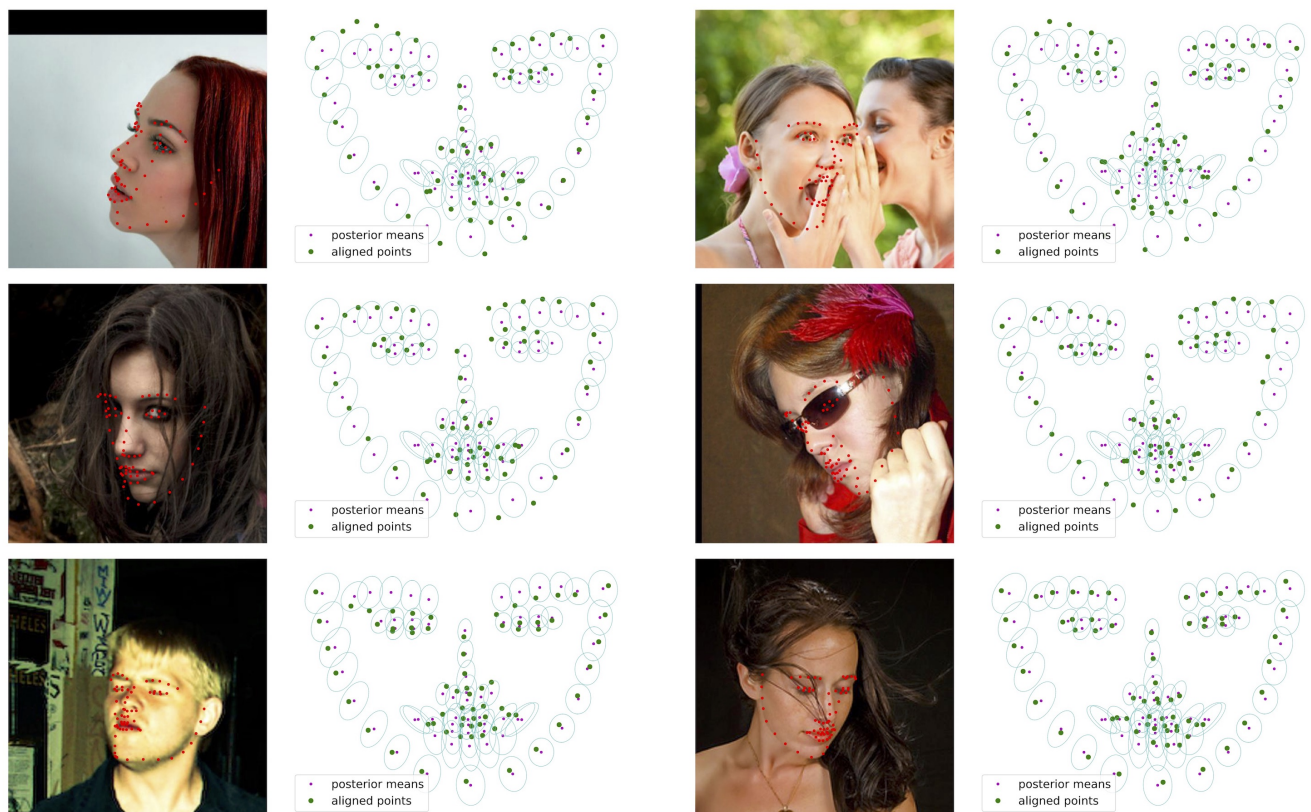


Fig. 15: Some examples of the worse scores obtained with GStudent-EM and the ground-truth landmarks.

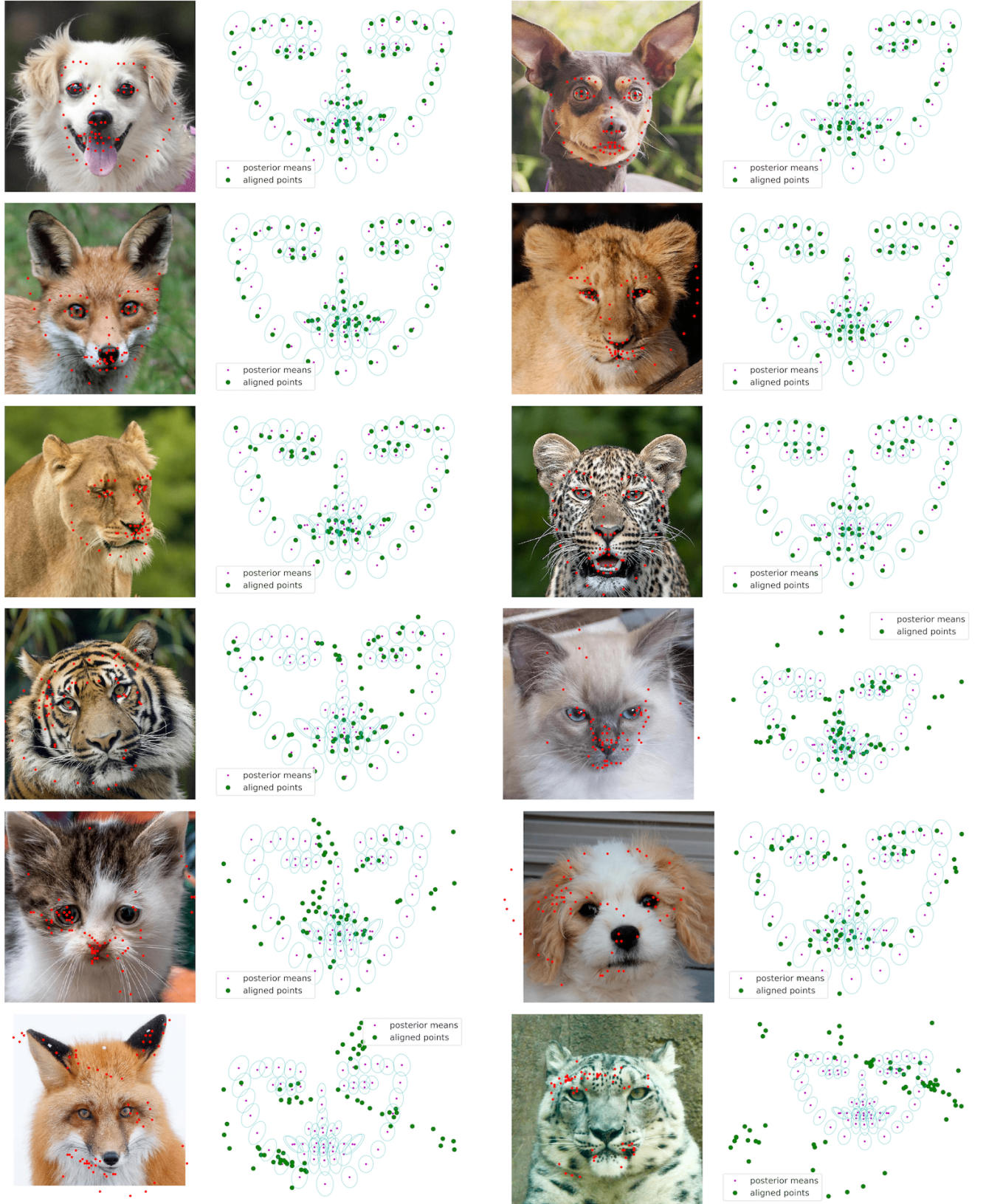


Fig. 16: Percentage of inliers obtained with 12 animal faces using 3DFA1 (Bulat and Tzimiropoulos, 2016) (from top to bottom and left to right: 0.91, 0.88, 0.88, 0.87, 0.75, 0.71, 0.13, 0.07, 0.06, 0.04, 0.0, 0.0).

Appendix A. Closed-Form Solution Using Unit Quaternions

Consider (17) with $\Sigma = \sigma \mathbf{I}$. We immediately obtain the following formulas for the model parameters:

$$s^* = \left(\frac{\sum_{n=1}^N \alpha_n \hat{\mathbf{y}}_n^\top \hat{\mathbf{y}}_n}{\sum_{n=1}^N \alpha_n \hat{\mathbf{x}}_n^\top \hat{\mathbf{x}}_n} \right)^{1/2}. \quad (\text{A.1})$$

$$\mathbf{R}^* = \underset{\mathbf{R}}{\operatorname{argmin}} \frac{1}{2} \sum_{n=1}^N \alpha_n \|\hat{\mathbf{y}}_n - s^* \mathbf{R} \hat{\mathbf{x}}_n\|^2, \quad (\text{A.2})$$

$$\sigma^* = \frac{1}{3 \sum_{n=1}^N \alpha_n} \sum_{n=1}^N \alpha_n \|\hat{\mathbf{y}}_n - s^* \mathbf{R}^* \hat{\mathbf{x}}_n\|^2, \quad (\text{A.3})$$

The formula for the posteriors becomes:

$$\alpha_n = \frac{p(2\pi\sigma)^{-3/2} \exp(-\|\mathbf{x}_n\|^2/2\sigma)}{p(2\pi\sigma)^{-3/2} \exp(-\|\mathbf{x}_n\|^2/2\sigma) + (1-p)\gamma^{-1}} \quad (\text{A.4})$$

It is well known that a rotation matrix can be parameterized by a unit quaternion [Horn \(1987\)](#). Let \mathbf{R} be parameterized by its axis and angle of rotation, $\mathbf{n} = (n_1 \ n_2 \ n_3)^\top$, $\|\mathbf{n}\| = 1$ and ϕ . The unit quaternion parameterizing the rotation is:

$$\begin{aligned} q &= \cos \frac{\phi}{2} + \sin \frac{\phi}{2} (in_1 + jn_2 + kn_3) \\ &= q_0 + iq_1 + jq_2 + kq_3, \end{aligned} \quad (\text{A.5})$$

with $i^2 = j^2 = k^2 = ijk = -1$, $\mathbf{q} = (q_0 \ q_1 \ q_2 \ q_3)^\top \in \mathbb{R}^4$ by abuse of notation, and $\mathbf{q}\mathbf{q}^\top = 1$. A vector $\mathbf{a} \in \mathbb{R}^3$ can be represented as a purely imaginary quaternion, namely $\tilde{\mathbf{a}} = (0 \ a_1 \ a_2 \ a_3)^\top \in \mathbb{R}^4$. The action of a rotation onto $\tilde{\mathbf{a}}$ can be written as $\mathbf{q} * \tilde{\mathbf{a}} * \bar{\mathbf{q}}$, where the symbol $*$ corresponds to the quaternion product and $\bar{\mathbf{q}}$ is the conjugate of \mathbf{q} , namely $\bar{\mathbf{q}} = q_0 - iq_1 - jq_2 - kq_3$. Making use of the properties $\|\mathbf{q}_1 * \mathbf{q}_2\|^2 = \|\mathbf{q}_1\|^2 \|\mathbf{q}_2\|^2$ and $\bar{\mathbf{q}} * \mathbf{q} = \|\mathbf{q}\|^2 = 1$, the squared Euclidean norm in (A.2) can be successively written as:

$$\begin{aligned} \|\hat{\mathbf{y}}_n - s\mathbf{R}\hat{\mathbf{x}}_n\|^2 &= \|\tilde{\mathbf{y}}_n - s\mathbf{q} * \tilde{\mathbf{x}}_n * \bar{\mathbf{q}}\|^2 \|\mathbf{q}\|^2 \\ &= \|\tilde{\mathbf{y}}_n * \mathbf{q} - s\mathbf{q} * \tilde{\mathbf{x}}_n * \bar{\mathbf{q}} * \mathbf{q}\|^2 \\ &= \|\tilde{\mathbf{y}}_n * \mathbf{q} - s\mathbf{q} * \tilde{\mathbf{x}}_n\|^2 \\ &= \|\mathbf{Q}(\tilde{\mathbf{y}}_n)\mathbf{q} - s\mathbf{W}(\tilde{\mathbf{x}}_n)\mathbf{q}\|^2 \\ &= \mathbf{q}^\top \mathbf{M}_n \mathbf{q}, \end{aligned} \quad (\text{A.6})$$

with:

$$\mathbf{M}_n = (\mathbf{Q}(\tilde{\mathbf{y}}_n) - s\mathbf{W}(\tilde{\mathbf{x}}_n))^\top (\mathbf{Q}(\tilde{\mathbf{y}}_n) - s\mathbf{W}(\tilde{\mathbf{x}}_n)), \quad (\text{A.7})$$

and where we replaced the quaternion products $\tilde{\mathbf{a}} * \mathbf{q}$ and $\mathbf{q} * \tilde{\mathbf{a}}$ with matrix-vector products, with:

$$\mathbf{Q}(\tilde{\mathbf{a}}) = \begin{pmatrix} 0 & -a_1 & -a_2 & -a_3 \\ a_1 & 0 & -a_3 & a_2 \\ a_2 & a_3 & 0 & -a_1 \\ a_3 & -a_2 & a_1 & 0 \end{pmatrix} \quad (\text{A.8})$$

$$\mathbf{W}(\tilde{\mathbf{a}}) = \begin{pmatrix} 0 & -a_1 & -a_2 & -a_3 \\ a_1 & 0 & a_3 & -a_2 \\ a_2 & -a_3 & 0 & a_1 \\ a_3 & a_2 & -a_1 & 0 \end{pmatrix} \quad (\text{A.9})$$

Consequently, the right-hand side of (A.2) writes

$$\sum_{n=1}^N (\mathbf{q}^\top \alpha_n \mathbf{M}_n \mathbf{q}) = \mathbf{q}^\top \left(\sum_{n=1}^N \alpha_n \mathbf{M}_n \right) \mathbf{q} = \mathbf{q}^\top \mathbf{M} \mathbf{q},$$

where $\alpha_n \geq 0$ and $\mathbf{M}_n \in \mathbb{R}^{4 \times 4}$ is semi-definite positive symmetric, i.e. (A.7), hence so is \mathbf{M} . By constraining the minimizer to be a unit quaternion, we obtain the following minimization problem:

$$\min_{\mathbf{q}} Q(\mathbf{q}) = \min_{\mathbf{q}} (\mathbf{q}^\top \mathbf{M} \mathbf{q} + \lambda(1 - \mathbf{q}^\top \mathbf{q})). \quad (\text{A.10})$$

From $dQ/d\mathbf{q} = 0$ we obtain $\mathbf{M}\mathbf{q}^* = \lambda\mathbf{q}^*$ and by substitution in (A.10) we obtain $Q(\mathbf{q}^*) = \lambda$. Therefore, the minimization problem (A.10) is equivalent to estimating the smallest eigenvalue-eigenvector pair $(\lambda^*, \mathbf{q}^*)$ of \mathbf{M} .

Acknowledgments. This work has been funded by the EU H2020 project #871245 SPRING and by the Multidisciplinary Institute in Artificial Intelligence (MIAI) # ANR-19-P3IA-0003.

References

- Abtahi S, Omidyeganeh M, Shirmohammadi S, Hariri B (2014) Yawdd: A yawning detection dataset. In: Proceedings of the 5th ACM Multimedia Systems Conference, pp 24–28 [2, 11](#)
- Arun KS, Huang TS, Blostein SD (1987) Least-squares fitting of two 3-D point sets. IEEE Transactions on Pattern Analysis and Machine Intelligence 9(5):698–700 [4, 5, 7](#)
- Bagdanov AD, Del Bimbo A, Masi I (2011) The Florence 2D/3D hybrid face dataset. In: Joint ACM Workshop on Human Gesture and Behavior Understanding, pp 79–80 [3](#)
- Banfield JD, Raftery AE (1993) Model-based gaussian and non-gaussian clustering. Biometrics pp 803–821 [4](#)
- Blanz V, Vetter T (1999) A morphable model for the synthesis of 3D faces. In: ACM SIGGRAPH, vol 99, pp 187–194 [1](#)
- Bonnans JF, Gilbert JC, Lemaréchal C, Sagastizábal CA (2006) Numerical optimization: theoretical and practical aspects. Springer Science & Business Media [4, 7](#)
- Bulat A, Tzimiropoulos G (2016) Two-stage convolutional part heatmap regression for the 1st 3D face alignment in the wild (3DFAW) challenge. In: European Conference on Computer Vision Workshops, Springer, pp 616–624 [2, 3, 11, 12, 13, 15, 18](#)
- Bulat A, Tzimiropoulos G (2017) How far are we from solving the 2D & 3D face alignment problem? (and a dataset of 230,000 3D facial landmarks). In: IEEE International Conference on Computer Vision, pp 1021–1030 [3, 11, 12, 13](#)
- Deng J, Roussos A, Chrysos G, Ververas E, Kotsia I, Shen J, Zafeiriou S (2019) The Menpo benchmark for multi-pose 2D and 3D facial landmark localisation and tracking. International Journal of Computer Vision 127(6-7):599–624 [2, 3](#)
- Escalera S, Baro X, Guyon I, Escalante HJ, Tzimiropoulos G, Valstar M, Pantic M, Cohn J, Kanade T (2018) Special issue on the computational face. IEEE Transactions on Pattern Analysis and Machine Intelligence 40(11):2541–2545 [1](#)
- Faugeras OD, Hebert M (1986) The representation, recognition, and locating of 3-d objects. The International Journal of Robotics Research 5(3):27–52 [4, 5](#)
- Feng Y, Wu F, Shao X, Wang Y, Zhou X (2018) Joint 3D face reconstruction and dense alignment with position map regression network. In: European Conference on Computer Vision, pp 534–551 [2, 3, 11, 12, 15](#)
- Forbes F, Wraith D (2014) A new family of multivariate heavy-tailed distributions with variable marginal amounts of tailweight: application to robust clustering. Statistics and computing 24(6):971–984 [4](#)

- Ghiasi G, Fowlkes CC (2014) Occlusion Coherence: Localizing Occluded Faces with a Hierarchical Deformable Part Model. In: IEEE Conference on Computer Vision and Pattern Recognition, pp 1899–1906 [11](#)
- Gou C, Wu Y, Wang FY, Ji Q (2016) Shape augmented regression for 3D face alignment. In: European Conference on Computer Vision, pp 604–615 [2](#)
- Gross R, Matthews I, Cohn J, Kanade T, Baker S (2010) Multi-PIE. Image and Vision Computing 28(5):807–813 [3](#)
- Hendrycks D, Gimpel K (2017) A baseline for detecting misclassified and out-of-distribution examples in neural networks. In: International Conference on Learning Representations [4](#)
- Hendrycks D, Mazeika M, Dietterich TG (2019) Deep anomaly detection with outlier exposure. In: International Conference on Learning Representations [4](#)
- Horaud R, Forbes F, Yguel M, Dewaele G, Zhang J (2010) Rigid and articulated point registration with expectation conditional maximization. IEEE Transactions on Pattern Analysis and Machine Intelligence 33(3):587–602 [4](#)
- Horn BK (1987) Closed-form solution of absolute orientation using unit quaternions. Journal of the Optical Society of America A 4(4):629–642 [4](#), [5](#), [6](#), [7](#), [8](#), [19](#)
- Horn BK, Hilden HM, Negahdaripour S (1988) Closed-form solution of absolute orientation using orthonormal matrices. Journal of the Optical Society of America A 5(7):1127–1135 [4](#), [5](#), [7](#)
- Jeni LA, Tulyakov S, Yin L, Sebe N, Cohn JF (2016) The first 3D face alignment in the wild (3DFAW) challenge. In: European Conference on Computer Vision, Springer, pp 511–520 [2](#), [3](#)
- Jeni LA, Cohn JF, Kanade T (2017) Dense 3D face alignment from 2D video for real-time use. Image and Vision Computing 58:13–24 [3](#)
- Jourabloo A, Liu X (2017) Pose-invariant face alignment via CNN-based dense 3D model fitting. International Journal of Computer Vision 124(2):187–203 [3](#)
- Kraft D (1988) A software package for sequential quadratic programming. Tech. Rep. DFVLR-FB 88-28, DLR German Aerospace Center – Institute for Flight Mechanics, Koln, Germany [7](#)
- Lathuilière S, Mesejo P, Alameda-Pineda X, Horaud R (2018) DeepGUM: Learning deep robust regression with a Gaussian-uniform mixture model. In: European Conference on Computer Vision, pp 202–217 [4](#)
- Lee K, Lee H, Lee K, Shin J (2018) Training confidence-calibrated classifiers for detecting out-of-distribution samples. In: International Conference on Learning Representations [4](#)
- Liang S, Li Y, Srikant R (2018) Enhancing the reliability of out-of-distribution image detection in neural networks. In: International Conference on Learning Representations [4](#)
- Loy CC, Liu X, Kim TK, De la Torre F, Chellappa R (2019) Special issue on deep learning for face analysis. International Journal of Computer Vision 127(6):533–536 [1](#)
- McLachlan G, Peel D (2000) Robust mixture modelling using the t distribution. Statistics and Computing 10(4):339–348 [4](#)
- Sanyal S, Bolkart T, Feng H, Black MJ (2019) Learning to regress 3D face shape and expression from an image without 3D supervision. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 7763–7772 [3](#)
- Savage LJ (1972) The Foundations of Statistics. Dover [10](#)
- Sun J, Kabán A, Garibaldi JM (2010) Robust mixture clustering using Pearson type VII distribution. Pattern Recognition Letters 31(16):2447–2454 [4](#), [6](#)
- Tu X, Zhao J, Jiang Z, Luo Y, Xie M, Zhao Y, He L, Ma Z, Feng J (2020) 3d face reconstruction from a single image assisted by 2d face images in the wild. IEEE Transactions on Multimedia [11](#), [16](#)
- Umeyama S (1991) Least-squares estimation of transformation parameters between two point patterns. IEEE Transactions on Pattern Analysis and Machine Intelligence 13(4):376–380 [4](#), [5](#)
- Wu Y, Ji Q (2019) Facial landmark detection: A literature survey. International Journal of Computer Vision 127(2):115–142 [1](#)
- Yin L, Sun XCY, Worm T, Reale M (2008) A high-resolution 3D dynamic facial expression database. In: IEEE International Conference on Automatic Face and Gesture Recognition [3](#)
- Yu R, Saito S, Li H, Ceylan D, Li H (2017) Learning dense facial correspondences in unconstrained images. In: The IEEE International Conference on Computer Vision (ICCV) [3](#)
- Zaharescu A, Horaud R (2009) Robust factorization methods using a Gaussian/uniform mixture model. International Journal of Computer Vision 81(3):240 [4](#)
- Zhang X, Yin L, Cohn JF, Canavan S, Reale M, Horowitz A, Liu P, Girard JM (2014) BP4D-spontaneous: a high-resolution spontaneous 3D dynamic facial expression database. Image and Vision Computing 32(10):692–706 [3](#)
- Zhao R, Wang Y, Benitez-Quiroz CF, Liu Y, Martinez AM (2016) Fast and precise face alignment and 3D shape reconstruction from a single 2D image. In: European Conference on Computer Vision Workshops, Springer, pp 590–603 [3](#)
- Zhu X, Lei Z, Liu X, Shi H, Li SZ (2016) Face alignment across large poses: a 3D solution. In: IEEE Conference on Computer Vision and Pattern Recognition, pp 146–155 [2](#), [3](#), [11](#), [12](#), [13](#), [16](#)
- Zhu X, Liu X, Lei Z, Li SZ (2019) Face alignment in full pose range: A 3d total solution. IEEE Transactions on Pattern Analysis and Machine Intelligence 41(1):78–92 [11](#)