# EMOCA: Emotion Driven Monocular Face Capture and Animation

Radek Daněček
rdanecek@tue.mpg.de

Michael J. Black
black@tue.mpg.de

Timo Bolkart
tbolkart@tue.mpg.de

Max Planck Institute for Intelligent Systems, Tübingen, Germany

Figure 1. **EMOCA** regresses 3D faces from images with facial geometry that captures the original emotional content. Top row: images of people with challenging expressions. Middle row: coarse shape reconstruction. Bottom row: reconstruction with detailed displacements.

## Abstract

*As 3D facial avatars become more widely used for communication, it is critical that they faithfully convey emotion. Unfortunately, the best recent methods that regress parametric 3D face models from monocular images are unable to capture the full spectrum of facial expression, such as subtle or extreme emotions. We find the standard reconstruction metrics used for training (landmark reprojection error, photometric error, and face recognition loss) are insufficient to capture high-fidelity expressions. The result is facial geometries that do not match the emotional content of the input image. We address this with EMOCA (EMOtion Capture and Animation), by introducing a novel deep perceptual emotion consistency loss during training, which helps ensure that the reconstructed 3D expression matches the expression depicted in the input image. While EMOCA achieves 3D reconstruction errors that are on par with the current best methods, it significantly outperforms them in terms of the quality of the reconstructed expression and the perceived emotional content. We also directly regress levels of valence and arousal and classify basic expressions from the estimated 3D face parameters. On the task of in-the-wild emotion recognition, our purely geometric approach is on par with the best image-based methods, highlighting the value of 3D geometry in analyzing human behavior. The model and code are publicly available at https://emoca.is.tue.mpg.de.*

## 1. Introduction

Teaching computers to see humans and understand their behavior is a long-standing goal of computer vision. To accomplish this, computers need to understand how humans look, how they move, and what they feel. Faces and their emotional expressions provide an important source of information about a person's internal emotional state. To support automated analysis of emotional state, we capture a person's face, including its 3D shape, pose, and facial expression, given a single RGB image. To do so, we go beyond prior work to extract 3D geometry that carries rich emotional content. We focus on parametric methods (i.e., animatable and model-based) due to their wide applicability for 3D avatar creation [38], image synthesis [34, 83], video editing [43, 86] and face recognition [9, 70].

The field of 3D face reconstruction has rapidly advanced over the last two decades; see Egger et al. [23] for a review. Existing methods that estimate 3D face models struggle to capture facial expressions in detail and often produce 3D shapes that do not carry the emotional content of the input image. This has several causes. First, some 3D face models lack sufficient expressiveness to capture subtle or extreme expressions. Second, reconstruction metrics like landmark reprojection loss [8], photometric loss [10], face recognition loss [32], or multi-image consistency losses [75, 82], are either not affected by facial expressions, or require perfect image alignment to capture subtle cues. Subtle changes

in geometry, however, can result in large differences in the perceived emotion. We argue that, to recover 3D expression accurately, we need a new reconstruction metric that measures differences in expressions between the 3D reconstruction and the input image.

To that end, we describe EMOCA (EMOtion Capture and Animation), a neural network that learns an animatable face model from in-the-wild images without 3D supervision. The design of our method is inspired by advances in the field of facial emotion recognition, which has made tremendous progress to date on estimating affect (or emotion) from in-the-wild-images [52]. Specifically, we train a state-of-the-art emotion recognition model, and leverage this during training of EMOCA as supervision. EMOCA introduces a novel perceptual *emotion consistency loss* that encourages the similarity of emotional content between the input and rendered reconstruction.

While the new emotion consistency loss results in better reconstructed emotions, this alone is insufficient. Large image datasets used by previous 3D reconstruction methods, while containing a large number of subjects of diverse ethnicities, lack emotional expressivity [14, 17, 46, 94]. Large datasets with facial expressions, valence, and arousal in-the-wild, on the other hand, while rich in emotions, do not provide multiple images per subject in diverse conditions [7, 13, 48–50, 59, 93] and smaller datasets in controlled settings are not suitable for deep learning [56–58, 62, 80]. Multiple images of the same person, however, are required to train current state-of-the-art 3D face reconstruction methods [20, 28, 75]. To overcome this, EMOCA builds on top of DECA [28], a publicly available 3D face reconstruction framework that achieves state-of-the-art identity shape reconstruction accuracy [30, 75]. Specifically, we augment DECA's architecture with an additional trainable prediction branch for facial expression, while keeping other parts fixed. This enables us to only train the expression part of EMOCA on emotion-rich image data [59], which results in improved emotion reconstruction performance, while retaining DECA's identity face shape quality.

Once trained, EMOCA reconstructs a 3D face from a single image (Fig. 1), it significantly outperforms previous state-of-the-art methods in terms of the reconstructed expression quality, it preserves the state-of-the-art identity shape reconstruction accuracy, and the reconstructed face can be readily animated. Further, the expression parameters regressed by EMOCA convey sufficient information for in-the-wild emotion recognition, with on-par performance with the best image-based methods [88].

In summary, our main contributions are: 1) The first approach to reconstruct an animatable 3D face model from an in-the-wild image, that is capable of recovering facial expressions that convey the correct emotional state. 2) A novel perceptual emotion-consistency loss that rewards the accuracy of the reconstructed emotion. 3) The first 3D geometry-based framework for in-the-wild emotion recognition, with comparable performance to current state-of-the-art image-based methods. 4) The code and model are publicly available for research purposes at `https://emoca.is.tue.mpg.de`.

## 2. Related work

**Monocular face reconstruction:** Reconstructing 3D face shape from images has been studied extensively for more than two decades [23, 101]. Model-free approaches directly regress 3D meshes [19, 21, 29, 35, 42, 72, 76, 81, 95, 97, 99] or voxels [40] from an image, or optimize a Signed Distance Function (SDF) [63] to fit a face image. Most of these methods require explicit 3D supervision during training. While the output is model-free, acquiring the training data typically relies on a 3D face model (3D Morphable Model, or 3DMM). Thus their ability to reconstruct expressive faces may be limited by the 3DMM-based reconstruction used to generate the paired training data [19, 29, 35, 40, 42, 72, 95], the domain gap between 3DMM-based synthetic training data and real images [21, 76, 99], or the regularization towards a fixed 3DMM fitting result [16]. Instead, EMOCA is trained in a self-supervised fashion without any explicit 3D supervision, which enables it to capture less constrained expressions. Other self-supervised methods do not leverage face-domain-specific knowledge, which makes them applicable to general objects, but also limits the reconstruction quality [81, 97]. Unlike EMOCA, none of these model-free methods separate facial identity from facial expression, making them inappropriate for applications like expression re-targeting or animation.

Several works reconstruct the parameters of fixed statistical models like the Basel Face Model (BFM) [64], FaceWarehouse [12], or FLAME [53], or jointly learn a model and reconstruct faces from images [82, 84, 91]. Existing methods can be categorized into optimization-based [4, 5, 9, 10, 33, 47, 65, 71, 86, 92] and learning-based. The latter are trained fully supervised [15, 36, 44, 69, 89, 90, 100] or self-supervised with predicted 2D keypoints [20, 28, 54, 75, 78, 82, 84, 85, 98], 2D face contours [54], photometric constraints [20, 28, 32, 78, 82, 84, 85, 98], face recognition features [20, 28, 32, 78], multi-view constraints [78], or multi-image constraints [20, 28, 32, 75, 82]. Each supervision signal impacts the reconstructed 3D face in a unique way. Explicit 3D mesh or model parameter supervision induces a bias towards the method used to generate the pseudo-ground truth. Using face recognition features or leveraging multiple images of the same identity during training mainly impacts identity shape and appearance. Keypoint losses impact the facial geometry and image alignment (global transformation, identity and expression shape parameters), but predicted keypoints are sparse (commonly 51-68 points),

often inaccurate - especially for extreme expressions and head poses - and obtaining the optimal embedding of the corresponding keypoints on the model's surface is challenging. Photometric losses impact all model parameters (global transformation, identity and expression shape, appearance, and lighting), but, as with the keypoint losses, are strongly affected by misalignments between the predicted 3D face and the image. While using multi-view data during training has the potential to reconstruct more accurate 3D faces, there are no large datasets with a large number of identities and large diversity in expression, ethnicity, age, lighting conditions, etc. Consequently, while the field of monocular in-the-wild face capture has made tremendous progress, there are still limitations, particularly in the accuracy of the reconstructed expressions, which limit the emotions that can be perceived from the reconstructed 3D shapes. EMOCA instead learns to reconstruct expressive faces by combining emotion features that mainly propagate to the reconstructed expression, with a unique self-supervised framework that enables us to leverage a large dataset of diverse expressions.

**Emotion analysis from images:** Emotion analysis is a long-standing problem in computer vision and related fields (see [3, 11] for comprehensive surveys). Emotional states are commonly represented as discrete basic [24, 25] (e.g., Happiness, Surprise, ...) or compound expression categories [22] (e.g., happily surprised), continuous valence (positive-negative) and arousal (relaxed-intensive) values [73], or Facial Action Units (FACS) activations [26], where each action unit (AU) corresponds to a particular emotion-related facial muscle movement.

Early work on expression recognition extracts geometric features defining shape and location of face components [61,87], appearance features [27,77], or combinations of these [41, Chapter 19]. Over the last decade, the availability of large datasets for single-image expression analysis [7, 59] and audio-visual videos [48–50] shifted the focus from manually designed features to end-to-end trained models [52]. While early work like Wen and Huang [96] uses 3D non-rigid surface tracking to extract features for expression reconstruction, the majority of 3D-based methods focus on recognizing expressions from 3D scans [60, 74]. Among these, the most relevant to EMOCA, is [67] as they use 3DMM features to classify three expressions (obtained by fitting the 3DMM to the scans); most other methods use diverse 2D and 3D features extracted from the textured 3D scans.

Few 3DMM-based methods exist to recognize expressions from images. Bejaoui et al. [6] fit a 3DMM to images, while Chang et al. [15] and Koujan et al. [51] train a 3DMM parameter regressor, fully-supervised by parameters obtained by fitting a 3DMM to images and videos. From the 3DMM expression parameters, they then learn to classify different expressions. Most related to EMOCA, Shi et al. [79] use an expression recognition loss during training, but with the goal of obtaining a more discriminative latent representation. These methods focus on recognizing expressions, not improving 3D reconstruction. In contrast, EMOCA leverages recent advances in emotion recognition to reconstruct more expressive 3D faces.

## 3. Preliminaries

**Face model:** FLAME [53] is a statistical 3D head model with parameters for identity shape $\boldsymbol{\beta} \in \mathbb{R}^{|\boldsymbol{\beta}|}$, facial expression $\boldsymbol{\psi} \in \mathbb{R}^{|\boldsymbol{\psi}|}$, and pose parameters $\boldsymbol{\theta} \in \mathbb{R}^{3k+3}$ for rotations around $k = 4$ joints (neck, jaw, and eyeballs) and the global rotation. Given all parameters, FLAME outputs a mesh with $n_v = 5023$ vertices. Formally, FLAME is:

$$M(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}) \rightarrow (\mathbf{V}, \mathbf{F}), \tag{1}$$

with vertices $\mathbf{V} \in \mathbb{R}^{n_v \times 3}$ and $n_f = 9976$ faces $\mathbf{F} \in \mathbb{R}^{n_f \times 3}$. FLAME comes with an appearance model, converted from Basel Face Model's albedo space [64] to FLAME's UV layout [1]. Given parameters $\boldsymbol{\alpha} \in \mathbb{R}^{|\boldsymbol{\alpha}|}$, this model outputs a FLAME texture map $A(\boldsymbol{\alpha}) \in \mathbb{R}^{d \times d \times 3}$.

**Face reconstruction:** DECA [28] is a publicly available framework to reconstruct a detailed, animatable 3D face model from a single image. We follow DECA's notation for simplicity. Given an image $I$, the coarse encoder

$$E_c(I) \rightarrow (\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}, \boldsymbol{\alpha}, \mathbf{l}, \mathbf{c}) \tag{2}$$

outputs FLAME geometry parameters $\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}$, albedo $\boldsymbol{\alpha}$, Spherical Harmonics (SH) [66] lighting $\mathbf{l} \in \mathbb{R}^{27}$, and camera $\mathbf{c} \in \mathbb{R}^3$, which is the concatenation of isotropic scale $s \in \mathbb{R}$ and translation $\mathbf{t} \in \mathbb{R}^2$. The detail encoder

$$E_d(I) \rightarrow \boldsymbol{\delta} \tag{3}$$

encodes $I$ to a subject-specific detail vector $\boldsymbol{\delta} \in \mathbb{R}^{128}$. To reconstruct dynamic expression wrinkles, the detail decoder

$$F_d(\boldsymbol{\delta}, \boldsymbol{\psi}, \boldsymbol{\theta}_{jaw}) \rightarrow D \tag{4}$$

uses $\boldsymbol{\delta}$ to parametrize static person-specific details, and FLAME's expression $\boldsymbol{\psi}$ and jaw-pose parameters $\boldsymbol{\theta}_{jaw}$ to generate an expression-dependent detail UV displacement map $D \in \mathbb{R}^{d \times d \times 3}$.

Denoting the rendering function with $R$ [68], the coarse shape can be rendered to a 2D image as $R(M(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}), \boldsymbol{\alpha}, \mathbf{l}, \mathbf{c}) \rightarrow I_{Rc}$. To render the FLAME mesh, with expression-dependent details, to an image $I_{Rd}$, the $D$ are converted to a detailed normal map $N_d$, and provided as additional parameters to $R$; formally $R(M(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}), \boldsymbol{\alpha}, \mathbf{l}, \mathbf{c}, N_d) \rightarrow I_{Rd}$.

**Relative keypoint loss:** Given 2D face keypoints $\mathbf{k}_i \in \mathbb{R}^2$ and the corresponding keypoints on FLAME's mesh surface $M_i \in \mathbb{R}^3$, the relative keypoint loss [28] computes

offset vectors between pairs of 2D keypoints and between the corresponding pairs of projected model keypoints, and penalizes the difference. Formally, the loss computes as

$$L_{rk}^E = \sum_{(i,j) \in E} \|\mathbf{k}_i - \mathbf{k}_j - s\Pi(M_i - M_j)\|_1, \quad (5)$$

where $E$ is a set of landmark index pairs, and $\Pi \in \mathbb{R}^{2 \times 3}$ is the orthographic 3D-2D projection matrix.

**Emotion recognition:** For the emotion network, we use ResNet-50 as the backbone, with a fully connected prediction head that outputs expression classification, valence, and arousal. See Appendix for experiments with other backbones. The network is trained on AffectNet [59], a large-scale annotated emotion dataset. We adapt the training setting from Toisoul et al. [88] with minor modifications as described in the Appendix. The loss function consists of several terms such as categorical cross entropy for expression classification, mean squared error and correlation coefficient losses for valence and arousal; see Appendix for details of the losses. After the network is trained, prediction heads are discarded, and the features of the final layer of the backbone network serve as our emotion feature $\boldsymbol{\epsilon} \in \mathbb{R}^{|\boldsymbol{\epsilon}|}$. We denote the emotion network as $A(I) \to \boldsymbol{\epsilon}$.

## 4. Method: EMOCA

The main goal of EMOCA is to address a significant limitation of the prior art - to recover 3D face shapes from single images that convey the full spectrum of emotion. Our technical contribution is twofold, first, we introduce a novel *emotion consistency loss* that is designed to encourage *emotion similarity* between the input image and the output rendering as training supervision. Second, we leverage parts of DECA's [28] trained model in order to only train the expression part of EMOCA on emotion-rich image data, while preserving DECA's identity shape reconstruction performance.

**Architecture:** EMOCA's architecture is based on DECA [28]. As with many state-of-the-art methods, DECA takes an input image and uses several neural networks to factor it into shape, albedo, lighting, etc. Given these factors, one can differentiably render an output image that should look like the input. Here we exploit this output image in a novel way by encouraging it to have the same *expression* as the input image.

Training models like DECA [28] on emotion-rich image data [59] is infeasible, due to DECA's requirement of multiple training images per subject to regularize the training of the identity shape reconstruction of $E_c$ (Eq. 2). Instead, we augment DECA's architecture with an additional expression encoder

$$E_e(I) \to \boldsymbol{\psi}_e, \quad (6)$$

and keep the weights of $E_c$ fixed during training, thereby retaining the predictions of $\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\alpha}, \mathbf{l}$ and $\mathbf{c}$ from DECA, but discarding DECA's $\boldsymbol{\psi}$. Further, let $R(M(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}_e), \boldsymbol{\alpha}, \mathbf{l}, \mathbf{c}) \to I_{Re}$ denote the rendering of the output of $E_c$ with the expression of the input image, $E_e(I)$.

For an overview of the model architecture, see Figure 2. Training $E_e$ only has several advantages: 1) Training datasets are not required to contain multiple images per subject. 2) Not training the identity prediction enables us to remove the face recognition loss. 3) Having fixed pose, shape, and camera parameters allows us to remove the landmark reprojection loss. 4) This results in reduced training resources, faster training time, and reduced memory consumption due to the lower number of training parameters.

**Loss function:** In total, we optimize:

$$\begin{aligned} L = & \lambda_{emo}L_{emo} + \lambda_{pho}L_{pho} + \lambda_{eye}L_{eye} \\ & + \lambda_{mc}L_{mc} + \lambda_{lc}L_{lc} + \lambda_{\boldsymbol{\psi}}L_{\boldsymbol{\psi}} \end{aligned} \quad (7)$$

with emotion consistency loss $L_{emo}$, photometric loss $L_{pho}$, eye closure loss $L_{eye}$, mouth closure loss $L_{mc}$, lip corner loss $L_{lc}$, and expression regularizer $L_{\boldsymbol{\psi}}$, each weighted by a factor $\lambda_x$.

**Emotion consistency loss:** The emotion consistency loss computes the difference between the emotion features of the input image $\boldsymbol{\epsilon}_I = A(I)$ and those of the rendered image, $\boldsymbol{\epsilon}_{Re} = A(I_{Re})$ as:

$$L_{emo} = d(\boldsymbol{\epsilon}_I, \boldsymbol{\epsilon}_{Re}), \quad (8)$$

with $d(\boldsymbol{\epsilon}_1, \boldsymbol{\epsilon}_2) = \|\boldsymbol{\epsilon}_1 - \boldsymbol{\epsilon}_2\|_2$. Instead of measuring a geometric error, $L_{emo}$ computes a perceptual difference between the input image and the rendered image. Optimizing this loss during training ensures that the reconstructed 3D face conveys the emotional content of the input image.

**Photometric loss:** The photometric loss computes the pixel error between the input image $I$ and the output rendering $I_{Re}$. $L_{pho} = \|V_I \odot (I - I_{Re})\|_{1,1}$. $V_I$ denotes a rendered mask of the output face shape, with each pixel located in the face skin region is equal to $1$, and $0$ elsewhere. The operator $\odot$ denotes the Hadamard product.

**Eye closure loss:** The eye closure loss computes as $L_{eye} = L_{rk}^{E_{eye}}$, where $E_{eye}$ is a set of upper/lower eyelid keypoint pairs. Due to slight misalignment between image landmarks and projected 3D landmarks, enforcing standard landmark reprojection losses produces incorrect predictions. Instead, using (translation-invariant) relative keypoint losses (for eye closure, mouth closure, and mouth width) is less susceptible to misalignments.

**Mouth closure loss:** The loss computes as $L_{mc} = L_{rk}^{E_{mc}}$, where $E_{mc}$ is a set of upper/lower lip keypoint pairs.

**Lip corner loss:** The lip corner loss computes as $L_{lc} = L_{rk}^{E_{lc}}$, where $E_{lc}$ is the pair of left and right lip corners.

**Expression regularization:** The expression is regularized as $L_{\boldsymbol{\psi}} = \|\boldsymbol{\psi}\|_2^2$.
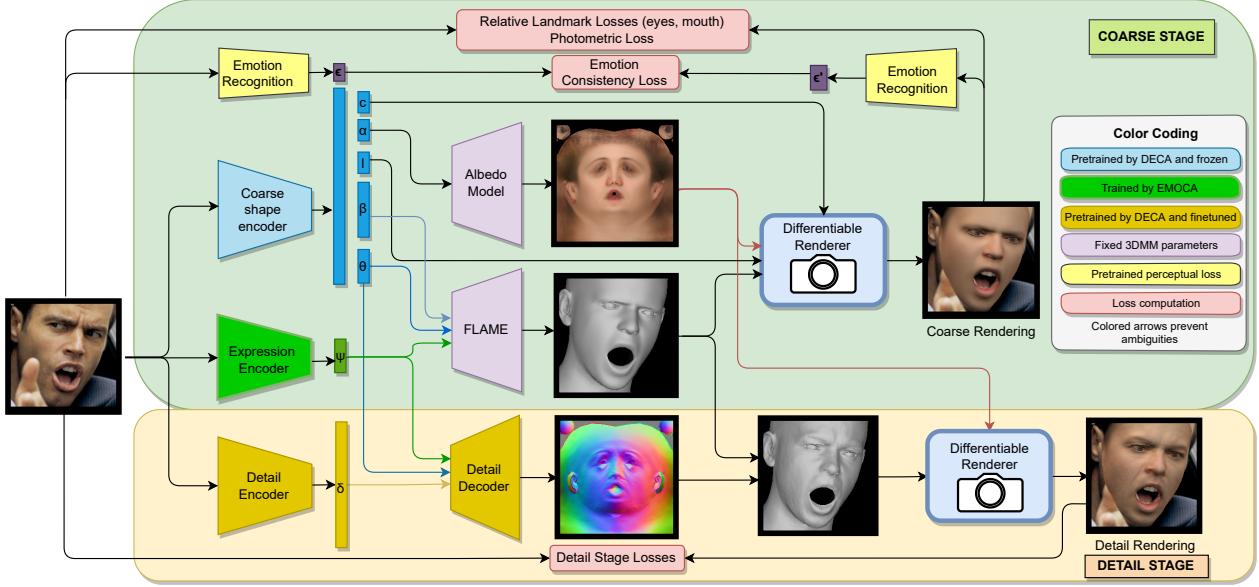
Figure 2. EMOCA overview. For the coarse training stage (green box), the input image is fed to the coarse shape encoder (initialized from DECA [28] and fixed) and EMOCA's trainable expression shape encoder. A textured 3D mesh is then reconstructed from the regressed identity shape, expression shape, pose, and albedo parameters with FLAME's geometry and albedo models as fixed decoders. This textured mesh is rendered by a differentiable renderer with the regressed camera and spherical harmonics lighting. Our novel emotion consistency loss (Eq. 8) penalizes the difference between the emotion features of the input image and those of the rendered coarse shape, after passing both images through a fixed emotion recognition network. For the detail training stage (yellow box), EMOCA's expression encoder is fixed, and the regressed expression (and jaw-pose) parameters are used to condition the detail decoder.

**Detailed stage:** The detail training stage adds wrinkle details that are animatable. Here we follow DECA's design, and use the same architecture and losses.

## 5. Experiments

### 5.1. Training setting

The first stage (coarse part) of EMOCA is trained with AffectNet [59] for a maximum of 20 epochs, with early stopping, using the Adam optimizer [45] and a learning rate of $5e - 5$. We use the same training/validation/testing split as proposed by [88]. We set $\lambda_{emo} = 1$, $\lambda_{pho} = 2$, $\lambda_{eye} = \lambda_{lc} = \lambda_{mc} = 0.5$ and $\lambda_{\psi} = 1e - 4$. EMOCA's second stage (detail part) training is comparable to DECA's second stage training. We use the same training data [14,17] and train with the same settings. Please refer to the Appendix for more training details.

### 5.2. Quantitative evaluation

While, for the task of 3D face reconstruction, standard benchmarks exist to quantitatively evaluate the identity face shape [30, 75], no such benchmark exists to assess the accuracy of the reconstructed expression. Unlike the identity shape benchmarks, quantitatively measuring the difference between a reconstructed 3D facial expression and a ground truth scan is less meaningful. The errors would be dominated by errors of the reconstructed identity face shape, and a low geometric error would not necessarily correspond to a small difference in human perception of the emotion. Instead, we evaluate EMOCA 1) qualitatively, 2) quantitatively for the task of in-the-wild emotion recognition, and 3) perceptually in an Amazon Mechanical Turk (AMT) study.

**Emotion recognition:** Our goal is to quantify how much of the input emotion is conveyed in the reconstructed 3D face. For this, we apply 3D face reconstruction methods to in-the-wild emotional face images and evaluate emotion recognition accuracy based on the 3D reconstruction. Here we focus on methods that reconstruct a parametric model of the face, i.e. a 3DMM. To that end, for each 3D face reconstruction method, we train a 4-layer MLP with Batch Normalization [39] and LeakyReLUs to regress valence and arousal levels, and classify expression labels directly from the predicted 3DMM parameters. The training details are described in the Appendix.

We evaluate emotion recognition on the AffectNet test set [59] and the AFEW-VA test set [49]. For each method, we report Concordance correlation coefficients (CCC ↑), Pearson correlation coefficients (PCC ↑), root mean squared error (RMSE ↓), and sign agreement (SAGR ↑) for valence (V) and arousal (A) regression and accuracy for expression (E) classification on the test set defined by [88]. EMOCA outperforms all 3D face reconstruction methods, and is on

| Model | V-PCC ↑ | V-CCC ↑ | V-RMSE ↓ | V-SAGR ↑ | A-PCC ↑ | A-CCC ↑ | A-RMSE ↓ | A-SAGR ↑ | E-ACC ↑ |
|---|---|---|---|---|---|---|---|---|---|
| EmoNet [88] | 0.75 | 0.73 | *0.32* | *0.80* | 0.68 | 0.65 | **0.29** | *0.78* | *0.68* |
| Deep3DFace [20] | 0.75 | 0.73 | 0.33 | *0.80* | 0.66 | 0.65 | 0.31 | 0.78 | 0.65 |
| ExpNet [15] | 0.45 | 0.42 | 0.43 | 0.73 | 0.39 | 0.36 | 0.38 | 0.64 | 0.46 |
| MGCNet [78] | 0.71 | 0.69 | 0.35 | *0.80* | 0.59 | 0.58 | 0.34 | 0.77 | 0.60 |
| 3DDFA_V2 [36] | 0.63 | 0.62 | 0.39 | 0.75 | 0.53 | 0.50 | 0.34 | 0.73 | 0.52 |
| DECA [28] | 0.70 | 0.69 | 0.36 | 0.76 | 0.59 | 0.58 | 0.33 | 0.74 | 0.59 |
| DECA w/ details [28] | 0.70 | 0.69 | 0.37 | 0.77 | 0.59 | 0.57 | 0.33 | 0.77 | 0.58 |
| EMOCA (Ours) | **0.78** | **0.77** | **0.31** | **0.81** | *0.69* | *0.68* | *0.30* | *0.81* | *0.68* |
| EMOCA w/ details (Ours) | *0.77* | *0.76* | **0.31** | **0.81** | **0.70** | **0.69** | **0.29** | **0.83** | **0.69** |

Table 1. **Emotion recognition performance on the AffectNet test set [59].** The EmoNet performance is measured using the model that is publicly released by the authors. For EMOCA and the other 3D baselines, we train the recognition module as described in Sec. 5.2. DECA w/ detail means that DECA's detail code prediction was included in the input to the regressor, along with the 3DMM parameters. Please note that EMOCA's performance is on par with EmoNet and it outperforms all other 3D reconstruction-based methods.
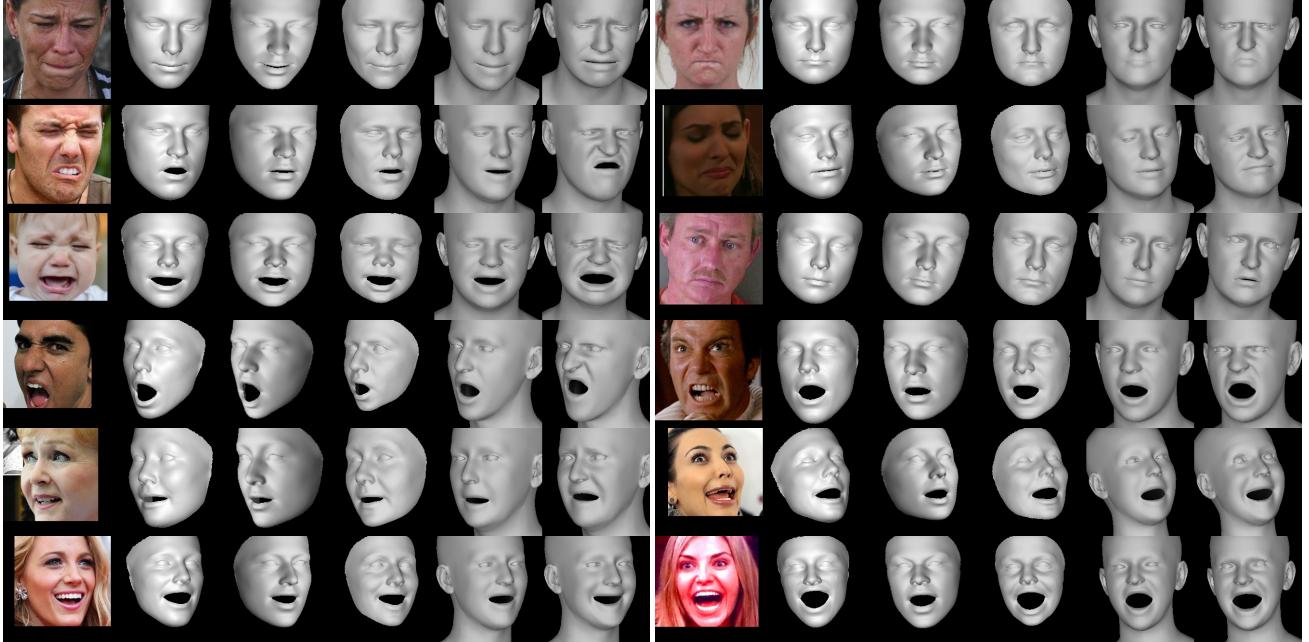


Figure 3. Comparison of **coarse reconstruction** methods, from left to right: Input, 3DDFA_V2 [36], MGCNet [78], Deng et al. [20], DECA [28] (coarse), and EMOCA (coarse). EMOCA conveys the emotions of the input images better than other methods.

par with the image-based state-of-the-art [88]. For details, see Tab. 1 for results on the AffectNet dataset and Tab. 3 of the Appendix for the AFEW-VA dataset.

Note that EMOCA performs on par with EmoNet [88], which is a recent method for estimating emotion from images. This confirms that the emotional content is present in our 3D reconstruction and that 3D shape is sufficient to understand emotion. This has implications for future research on emotion recognition.

**Perceptual study:** The 3D geometry reconstructed from an image must convey the emotion of the input image. Directly comparing rendered geometry with an image is difficult due to the domain gap. Instead, we perform a perceptual study using AMT to assess the perceived expression of emotion from rendered 3D reconstructions. Specifically,

given an image, we ask participants to categorize the perceived expression of emotion into one of the 7 basic emotions (Anger, Disgust, Fear, Happiness, Sadness, Surprise, and Contempt) or as a neutral expression (no emotion). A single evaluation task contains 75 images in random order; 35 real images, the 35 corresponding rendered reconstructions (from one method), and 5 qualification samples. The 5 qualification samples are duplicates sampled from the 35 real images, and they are chosen to be of easily recognizable emotion. Each task is performed by 10 participants. Participants that either misclassify the emotion, or inconsistently label the duplicated images, for at least 2 (of the 5) qualification samples are discarded from further analysis to filter out inattentive and/or uncooperative participants. For each method, we measure the classification consistency between

each participant's labels for the rendered images and their labels for the corresponding real images. If the rendered 3D meshes contain the emotional content of the images, then the scores given to both should be consistent.

We select 35 images with balanced emotional content (i.e., 5 images per basic emotion) from the AffectNet test set [59]. For each image, we reconstruct 3D faces using EMOCA, DECA [28], Deep3DFace [20], MGCNet [78] and 3DDFA_V2 [36]. The classification consistency averaged across participants for each method are: EMOCA (coarse) 0.68, EMOCA (detail) 0.65, Deep3DFace 0.37, DECA (coarse) 0.33, DECA (detail) 0.31, MGCNet 0.32, 3DDFA_V2 0.31. In summary, EMOCA preserves the emotional content of images better than the other methods. Note that, perhaps surprisingly, there is little difference between the scores for the coarse meshes of EMOCA/DECA and the detailed ones. Despite more wrinkle detail, our perceptual experiments suggest that the detailed meshes do not convey more emotional content. One possible explanation is that, in addition to adding valid wrinkle details, the detail generator sometimes adds artifacts in the lip region (e.g. Fig. 1, col. 1 & 3), and hallucinates details in the forehead (e.g. Fig. 1, col. 8). These could negatively impact participants' perception. For the full confusion tables, see the Appendix.

**Emotion recognition vs. perceptual study:** There is a considerable discrepancy between the results of the automatic emotion recognition results and the perceptual study results, in particular for Deep3DFace [20]. Deep3DFace performs much better on the emotion recognition task (slightly below SOTA), than on the perceptual study. Unlike EMOCA, it is not capable of producing highly emotional reconstructions (see Fig. 3). We hypothesize that the automatic predictors are capable of detecting more subtle cues than humans. We investigate this by measuring the agreement (i.e., percentage of matching predictions) between the method's classifier (from the reconstructed face parameters) and the participant's annotation of the *input images* from the perceptual study. The results are: EMOCA 62% and Deep3DFace 62%. This indicates that the predicted parameters for both methods contain a similar amount of information about the emotions compared to the annotations of the input images. However, the agreement between the method's classifier and the participant's annotation of the *rendered reconstructions* is for EMOCA 48%, and for Deep3DFace 26%. In other words, EMOCA is signficantly more in agreement with human perception.

### 5.3. Qualitative evaluation

We provide a visual comparison of the coarse shape reconstruction methods in Fig. 3. Observe that EMOCA outperforms all the previous methods in terms of capturing the emotional content of the original image in the reconstructed expression. In Fig. 4 we compare our detail reconstruc-



Figure 4. Comparison of 3D reconstructions with **detail displacements**. Top: Input, Middle: DECA [28], Bottom: EMOCA. EMOCA results contain more expression-dependent details that better convey the emotion of the input images than DECA.

tions to DECA's detail reconstruction. Compared to DECA, our detailed displacement better captures the fine details of highly emotional input images.

### 5.4. Ablation experiment

Table 2 shows the effect of ablating the training data and the emotion consistency loss. The table summarizes the effect of EMOCA trained w/ and w/o the emotion consistency loss, and using the DECA data only [14, 17] instead of the AffectNet training data [59].

## 6. Discussion and limitations

**Baseline:** EMOCA builds on top of DECA due to its state-of-the-art identity shape reconstruction performance. We found in our experiments that the recently released Deep3DFaceRecon [2] gives better 3D face reconstructions than reported in the paper [20], and in some cases, it outperforms DECA in terms of the reconstructed expression. Combining our emotion consistency loss with the Deep3DFaceRecon framework to further improve their reconstructed expressions is worth further investigation.

**Image alignment:** DECA sometimes predicts 3D faces that are slightly misaligned with the input images. EMOCA inherits this limitation due to the fixed coarse shape encoder. Further, while EMOCA reconstructs more expressive faces that better convey the emotion of the input image, expressions are also sometimes misaligned. Mitigating these artifacts, and better balancing the trade-off between geometric alignment and emotion similarity, requires further work.

**Emotion embedding analysis:** We assume that the emotion embedding extracted by the emotion recognition network has desirable properties to guide the optimization of FLAME's expression parameters. We found that the emotion recognition loss is more difficult to optimize and it requires more careful weighting of the loss compared to the identity recognition losses used in previous work [20, 28, 32]. Directly using the pre-trained EmoNet [88]

| Model | V-PCC ↑ | V-CCC ↑ | V-RMSE ↓ | V-SAGR ↑ | A-PCC ↑ | A-CCC ↑ | A-RMSE ↓ | A-SAGR ↑ | E-ACC ↑ |
|---|---|---|---|---|---|---|---|---|---|
| DECA [28] | 0.70 | 0.69 | 0.36 | 0.76 | 0.59 | 0.58 | 0.33 | 0.74 | 0.59 |
| EMOCA DS w/o Emo | 0.70 | 0.69 | 0.37 | 0.78 | 0.61 | 0.58 | 0.32 | *0.79* | 0.60 |
| EMOCA w/o Emo | 0.68 | 0.66 | 0.36 | 0.74 | 0.59 | 0.58 | 0.32 | 0.77 | 0.59 |
| EMOCA DS | *0.77* | *0.76* | **0.31** | **0.82** | **0.69** | *0.67* | **0.29** | *0.79* | **0.68** |
| EMOCA | **0.78** | **0.77** | **0.31** | *0.81* | **0.69** | **0.68** | *0.30* | **0.81** | **0.68** |

Table 2. **Ablation experiment.** Effect of ablating the data and the emotion consistency loss on EMOCA evaluated on the emotion recognition task. From top to bottom, we see the performance of DECA, EMOCA trained on the DECA dataset w/o emotion loss, EMOCA w/o emotion loss, EMOCA trained on the DECA dataset, and EMOCA. We refer to DECA's training data as DECA dataset (DS), which is a combination of VGGFace2 [14], and VoxCeleb2 [17]. The key finding is that novel emotion consistency loss is critical for the performance of this task, as with emotion loss, EMOCA's performance improves. Finetuning on AffectNet, which has a much richer variety of facial expressions, only marginally increases in performance over training on DECA's original training data (DS).

for instance did not provide sufficient supervision. However, our work is the first to demonstrate how to use emotion recognition features to guide the task of 3D geometry reconstruction. In addition using our emotion consistency loss to train EMOCA, we have experimented with the applicability of emotion features for the tasks of emotion retrieval and emotion retargetting via FLAME expression parameter optimization (see Appendix).

**Emotion network architecture:** Using a pre-trained state-of-the-art emotion recognition network [88] does not provide satisfactory supervision during optimization or training. Instead, it produces strong artifacts in the reconstructed geometry. To overcome this, we investigate different ResNet [37] and Swin Transformer [55] based emotion network architectures, and show the effect of different networks in the Appendix. Based on this analysis, we use a ResNet-50 backbone for our emotion network.

**Jaw rotations:** While FLAME's jaw rotation parameters $\theta_{jaw}$ contribute to facial expressions, we found the optimization of $\theta_{jaw}$ to be unstable while training EMOCA. We hypothesize, that this is due to the lack of a good prior for the jaw rotation. However, using different simplified priors for the jaw pose like a simple L2 regularizer did not give satisfactory results. We offer a more detailed discussion in the Appendix. Investigating the effect of more advanced data-driven jaw priors when optimizing the emotion loss is subject to future work.

**Implementation details:** For details on all hyper parameters and discussion on design choices see the Appendix.

## 7. Conclusions

We have presented EMOCA, a method that takes a single in-the-wild image and reconstructs a 3D face with sufficient facial expression detail to convey the emotional state of the input image. EMOCA is trained in a self-supervised fashion from a large dataset of emotion-rich images. A novel *emotion similarity* loss provides supervision on the reconstructed expressions during training. The emotion similarity relies on deep features extracted from a neural network trained for single-image affect (emotion) recognition in-the-wild. EMOCA reconstructs 3D face shape on par with current state-of-the-art methods but outperforms them in terms of the quality of the reconstructed expression. Further, using the reconstructed expression parameters for the task of in-the-wild emotion recognition, EMOCA outperforms existing 3DMM-based face reconstruction methods and gives on par results with the best purely image-based method.

In summary, this is the first in-the-wild monocular face reconstruction work that puts explicit emphasis on the *perceptual quality* of the expression and the emotion it communicates instead of standard geometric and photometric losses. This presents a new direction for the monocular face reconstruction community. This work has potential to further combine the fields of monocular 3D face reconstruction and emotion analysis. Further, downstream application of this work can be employed in the industry, including but not limited to gaming, movies, AR/VR and communication.

Of course, any improvement to 3D face acquisition and animation may also enable more realistic 'deep fakes.' Subtle emotional cues are individualistic and reproducing these could make it harder to detect such fakes. While cognizant of the risks, we are also sensitive the the importance of facial emotion in human communication. The trend towards emotional avatars in games and communication is clear. If communicative avatars do not properly communicate emotion, that, in itself presents a risk of misunderstandings.

# References

[1] BFM_to_FLAME. https://github.com/TimoBolkart/BFM_to_FLAME, 2021. 3

[2] Deep3DFaceRecon_PyTorch. https://github.com/microsoft/Deep3DFaceReconstruction, 2021. 7

[3] Xavier Alameda-Pineda, Elisa Ricci, and Nicu Sebe. Multimodal behavior analysis in the wild: An introduction. In Xavier Alameda-Pineda, Elisa Ricci, and Nicu Sebe, editors, *Multimodal Behavior Analysis in the Wild*, Computer Vision and Pattern Recognition, pages 1–8. Academic Press, 2019. 3

[4] Oswald Aldrian and William AP Smith. Inverse rendering of faces with a 3D morphable model. *Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 35(5):1080–1093, 2013. 2

[5] Anil Bas, William A. P. Smith, Timo Bolkart, and Stefanie Wuhrer. Fitting a 3D morphable model to edges: A comparison between hard and soft correspondences. In *Asian Conference on Computer Vision Workshops*, pages 377–391, 2017. 2

[6] Hela Bejaoui, Haythem Ghazouani, and Walid Barhoumi. Fully automated facial expression recognition using 3D morphable model and mesh-local binary pattern. In *Advanced Concepts for Intelligent Vision Systems*, volume 10617, pages 39–50, 2017. 3

[7] Carlos Fabian Benitez-Quiroz, Ramprakash Srinivasan, and Aleix M. Martínez. EmotioNet: An accurate, real-time algorithm for the automatic annotation of a million facial expressions in the wild. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5562–5570, 2016. 2, 3

[8] Volker Blanz, Curzio Basso, Tomaso A. Poggio, and Thomas Vetter. Reanimating faces in images and video. *Computer Graphics Forum (Proc. EUROGRAPHICS)*, 22(3):641–650, 2003. 1

[9] Volker Blanz, Sami Romdhani, and Thomas Vetter. Face identification across different poses and illuminations with a 3D morphable model. In *International Conference on Automatic Face & Gesture Recognition (FG)*, pages 202–207, 2002. 1, 2

[10] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3D faces. In *SIGGRAPH*, pages 187–194, 1999. 1, 2

[11] Rafael A. Calvo, Sidney D'Mello, Jonathan Gratch, and Arvid Kappas. *The Oxford Handbook of Affective Computing*. Oxford University Press, Inc., USA, 1st edition, 2014. 3

[12] Chen Cao, Yanlin Weng, Shun Zhou, Yiying Tong, and Keliang Zhou. FaceWarehouse: A 3D facial expression database for visual computing. *Transactions on Visualization and Computer Graphics*, 20:413–425, 2014. 2

[13] Houwei Cao, David Cooper, Michael Keutmann, Ruben Gur, Ani Nenkova, and Ragini Verma. Crema-d: Crowd-sourced emotional multimodal actors dataset. *IEEE transactions on affective computing*, 5:377–390, 10 2014. 2

[14] Qiong Cao, Li Shen, Weidi Xie, Omkar M Parkhi, and Andrew Zisserman. VGGFace2: A dataset for recognising faces across pose and age. In *International Conference on Automatic Face & Gesture Recognition (FG)*, pages 67–74, 2018. 2, 5, 7, 8, 14

[15] Feng-Ju Chang, Anh Tuan Tran, Tal Hassner, Iacopo Masi, Ram Nevatia, and Gerard Medioni. ExpNet: Landmark-free, deep, 3D facial expressions. In *International Conference on Automatic Face & Gesture Recognition (FG)*, pages 122–129, 2018. 2, 3, 6

[16] Aggelina Chatziagapi, ShahRukh Athar, Francesc Moreno-Noguer, and Dimitris Samaras. SIDER: single-image neural optimization for facial geometric detail recovery. In *International Conference on 3D Vision (3DV)*, pages 815–824, 2021. 2

[17] Joon Son Chung, Arsha Nagrani, and Andrew Zisserman. VoxCeleb2: Deep speaker recognition. In *Conference of the International Speech Communication Association (INTERSPEECH)*, pages 1086–1090, 2018. 2, 5, 7, 8, 14

[18] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, 2009. 14

[19] Jiankang Deng, Jia Guo, Evangelos Ververas, Irene Kotsia, and Stefanos Zafeiriou. RetinaFace: Single-shot multi-level face localisation in the wild. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5202–5211, 2020. 2

[20] Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. Accurate 3D face reconstruction with weakly-supervised learning: From single image to image set. In *Conference on Computer Vision and Pattern Recognition Workshops (CVPR-W)*, pages 285–295, 2019. 2, 6, 7, 16

[21] Pengfei Dou, Shishir K Shah, and Ioannis A Kakadiaris. End-to-end 3D face reconstruction with deep neural networks. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5908–5917, 2017. 2

[22] Shichuan Du, Yong Tao, and Aleix M Martinez. Compound facial expressions of emotion. *National Academy of Sciences*, 111(15):E1454–E1462, 2014. 3

[23] Bernhard Egger, William A. P. Smith, Ayush Tewari, Stefanie Wuhrer, Michael Zollhöfer, Thabo Beeler, Florian Bernard, Timo Bolkart, Adam Kortylewski, Sami Romdhani, Christian Theobalt, Volker Blanz, and Thomas Vetter. 3D morphable face models - past, present, and future. *Transactions on Graphics (TOG)*, 39(5):157:1–157:38, 2020. 1, 2

[24] Paul Ekman. An argument for basic emotions. *Cognition and Emotion*, 6(3-4):169–200, 1992. 3

[25] Paul Ekman and Wallace V Friesen. Constants across cultures in the face and emotion. *Journal of personality and social psychology*, 17(2):124, 1971. 3

[26] Paul Ekman and Wallace V. Friesen. Facial action coding system: A technique for the measurement of facial movement. *Palo Alto: Consulting Psychologists Press*, 1978. 3

9

[27] Xiaoyi Feng, M Pietikainen, and Abdenour Hadid. Facial expression recognition with local binary patterns and linear programming. *Pattern Recognition And Image Analysis*, 15(2):546, 2005. 3

[28] Yao Feng, Haiwen Feng, Michael J. Black, and Timo Bolkart. Learning an animatable detailed 3D face model from in-the-wild images. *Transactions on Graphics, (Proc. SIGGRAPH)*, 40(4):88:1–88:13, 2021. 2, 3, 4, 5, 6, 7, 8, 14

[29] Yao Feng, Fan Wu, Xiaohu Shao, Yanfeng Wang, and Xi Zhou. Joint 3D face reconstruction and dense alignment with position map regression network. In *European Conference on Computer Vision (ECCV)*, pages 534–551, 2018. 2

[30] Zhen-Hua Feng, Patrik Huber, Josef Kittler, Peter Hancock, Xiao-Jun Wu, Qijun Zhao, Paul Koppen, and Matthias Rätsch. Evaluation of dense 3D reconstruction from 2D face images in the wild. In *International Conference on Automatic Face & Gesture Recognition (FG)*, pages 780–786, 2018. 2, 5

[31] Xavier Gastaldi. Shake-shake regularization. *CoRR*, abs/1705.07485, 2017. 13

[32] Kyle Genova, Forrester Cole, Aaron Maschinot, Aaron Sarna, Daniel Vlasic, and William T. Freeman. Unsupervised training for 3D morphable model regression. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8377–8386, 2018. 1, 2, 7

[33] Thomas Gerig, Andreas Morel-Forster, Clemens Blumer, Bernhard Egger, Marcel Luthi, Sandro Schönborn, and Thomas Vetter. Morphable face models-an open framework. In *International Conference on Automatic Face & Gesture Recognition (FG)*, pages 75–82, 2018. 2

[34] Partha Ghosh, Pravir Singh Gupta, Roy Uziel, Anurag Ranjan, Michael J. Black, and Timo Bolkart. GIF: Generative interpretable faces. In *International Conference on 3D Vision (3DV)*, pages 868–878, 2020. 1

[35] Riza Alp Güler, George Trigeorgis, Epameinondas Antonakos, Patrick Snape, Stefanos Zafeiriou, and Iasonas Kokkinos. DenseReg: Fully convolutional dense shape regression in-the-wild. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6799–6808, 2017. 2

[36] Jianzhu Guo, Xiangyu Zhu, Yang Yang, Fan Yang, Zhen Lei, and Stan Z Li. Towards fast, accurate and stable 3D dense face alignment. In *European Conference on Computer Vision (ECCV)*, pages 152–168, 2020. 2, 6, 7, 16

[37] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 8, 14, 17, 21

[38] Liwen Hu, Shunsuke Saito, Lingyu Wei, Koki Nagano, Jaewoo Seo, Jens Fursund, Iman Sadeghi, Carrie Sun, Yen-Chun Chen, and Hao Li. Avatar digitization from a single image for real-time rendering. *Transactions on Graphics (TOG)*, 36(6):195:1–195:14, 2017. 1

[39] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning (ICML)*, volume 37, pages 448–456, 2015. 5

[40] Aaron S Jackson, Adrian Bulat, Vasileios Argyriou, and Georgios Tzimiropoulos. Large pose 3D face reconstruction from a single image via direct volumetric CNN regression. In *International Conference on Computer Vision (ICCV)*, pages 1031–1039, 2017. 2

[41] Anil K Jain and Stan Z Li. *Handbook of face recognition*, volume 1. Springer, 2011. 3

[42] Harim Jung, Myeong-Seok Oh, and Seong-Whan Lee. Learning free-form deformation for 3D face reconstruction from in-the-wild images. In *International Conference on Systems, Man, and Cybernetics (SMC)*, pages 2737–2742, 2021. 2

[43] Hyeongwoo Kim, Pablo Garrido, Ayush Tewari, Weipeng Xu, Justus Thies, Matthias Nießner, Patrick Pérez, Christian Richardt, Michael Zollhöfer, and Christian Theobalt. Deep video portraits. *Transactions on Graphics (TOG)*, 37(4):163:1–163:14, 2018. 1

[44] Hyeongwoo Kim, Michael Zollhöfer, Ayush Tewari, Justus Thies, Christian Richardt, and Christian Theobalt. InverseFaceNet: deep monocular inverse face rendering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4625–4634, 2018. 2

[45] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. 5

[46] Brendan F. Klare, Ben Klein, Emma Taborsky, Austin Blanton, Jordan Cheney, Kristen Allen, Patrick Grother, Alan Mah, Mark James Burge, and Anil K. Jain. Pushing the frontiers of unconstrained face detection and recognition: IARPA janus benchmark A. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1931–1939, 2015. 2

[47] Tatsuro Koizumi and William A. P. Smith. "look ma, no landmarks!" - unsupervised, model-based dense face alignment. In *European Conference on Computer Vision (ECCV)*, volume 12347, pages 690–706, 2020. 2

[48] Dimitrios Kollias and Stefanos Zafeiriou. Aff-Wild2: Extending the Aff-Wild database for affect recognition. *CoRR*, abs/1811.07770, 2018. 2, 3, 22, 26

[49] Jean Kossaifi, Georgios Tzimiropoulos, Sinisa Todorovic, and Maja Pantic. Afew-va database for valence and arousal estimation in-the-wild. *Image and Vision Computing*, 65:23–36, 2017. 2, 3, 5, 14, 16

[50] Jean Kossaifi, Robert Walecki, Yannis Panagakis, Jie Shen, Maximilian Schmitt, Fabien Ringeval, Jing Han, Vedhas Pandit, Antoine Toisoul, Björn W. Schuller, Kam Star, Elnar Hajiyev, and Maja Pantic. SEWA DB: A rich database for audio-visual emotion and sentiment research in the wild. *Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 43(3):1022–1040, 2021. 2, 3

[51] Mohammad Rami Koujan, Luma Alharbawee, Giorgos Giannakakis, Nicolas Pugeault, and Anastasios Roussos. Real-time facial expression recognition "in the wild" by disentangling 3d expression from identity. In *International Conference on Automatic Face & Gesture Recognition (FG)*, pages 24–31, 2020. 3

[52] Shan Li and Weihong Deng. Deep facial expression recognition: A survey. *Transactions on Affective Computing*, 2020. 2, 3

10

[53] Tianye Li, Timo Bolkart, Michael. J. Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4D scans. *Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 36(6):194:1–194:17, 2017. 2, 3, 17

[54] Yaojie Liu, Amin Jourabloo, William Ren, and Xiaoming Liu. Dense face alignment. In *International Conference on Computer Vision Workshops (ICCV-W)*, pages 1619–1628, 2017. 2

[55] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin Transformer: Hierarchical vision transformer using shifted windows. In *International Conference on Computer Vision (ICCV)*, pages 9992–10002, 2021. 8, 14, 17, 21, 22

[56] Patrick Lucey, Jeffrey F. Cohn, Takeo Kanade, Jason Saragih, Zara Ambadar, and Iain Matthews. The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression. In *Conference on Computer Vision and Pattern Recognition Workshops (CVPR-W)*, pages 94–101, 2010. 2

[57] Mohammad Mavadati, Peyten Sanger, and Mohammad H. Mahoor. Extended DISFA dataset: Investigating posed and spontaneous facial expressions. In *Conference on Computer Vision and Pattern Recognition Workshops (CVPR-W)*, pages 1452–1459, 2016. 2

[58] S. Mohammad Mavadati, Mohammad H. Mahoor, Kevin Bartlett, Philip Trinh, and Jeffrey F. Cohn. Disfa: A spontaneous facial action intensity database. *IEEE Transactions on Affective Computing*, 4(2):151–160, 2013. 2

[59] Ali Mollahosseini, Behzad Hasani, and Mohammad H Mahoor. Affectnet: A database for facial expression, valence, and arousal computing in the wild. *IEEE Transactions on Affective Computing*, 10(1):18–31, 2017. 2, 3, 4, 5, 6, 7, 14

[60] Francesca Nonis, Nicole Dagnes, Federica Marcolin, and Enrico Vezzetti. 3D approaches and challenges in facial expression recognition algorithms — a literature review. *Applied Sciences*, 9(18):3904, 2019. 3

[61] Maja Pantic and Léon J. M. Rothkrantz. Expert system for automatic analysis of facial expressions. *Image and Vision Computing*, 18(11):881–905, 2000. 3

[62] M. Pantic, Michel Valstar, R. Rademaker, and L. Maat. Web-based database for facial expression analysis. In *International Conference on Multimedia and Expo (ICME)*, pages 317–321, 2005. 2

[63] Jeong Joon Park, Peter Florence, Julian Straub, Richard A. Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 165–174, 2019. 2

[64] Pascal Paysan, Reinhard Knothe, Brian Amberg, Sami Romdhani, and Thomas Vetter. A 3D face model for pose and illumination invariant face recognition. In *International Conference on Advanced Video and Signal based Surveillance (AAAI)*, pages 296–301, 2009. 2, 3, 16

[65] Stylianos Ploumpis, Evangelos Ververas, Eimear O' Sullivan, Stylianos Moschoglou, Haoyang Wang, Nick E. Pears, William A. P. Smith, Baris Gecer, and Stefanos Zafeiriou. Towards a complete 3D morphable model of the human head. *Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 43(11):4142–4160, 2021. 2

[66] Ravi Ramamoorthi and Pat Hanrahan. An efficient representation for irradiance environment maps. *Annual Conference on Computer Graphics and Interactive Techniques*, pages 497–500, 2001. 3

[67] Subramanian Ramanathan, Ashraf A. Kassim, Y. V. Venkatesh, and Wu Sin Wah. Human facial expression recognition using a 3D morphable model. In *International Conference on Image Processing (ICIP)*, pages 661–664, 2006. 3

[68] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3D deep learning with pytorch3d. *arXiv:2007.08501*, 2020. 3

[69] E. Richardson, M. Sela, and R. Kimmel. 3D face reconstruction by learning from synthetic data. In *International Conference on 3D Vision (3DV)*, pages 460–469, 2016. 2

[70] Sami Romdhani, Volker Blanz, and Thomas Vetter. Face identification by fitting a 3D morphable model using linear shape and texture error functions. In *European Conference on Computer Vision (ECCV)*, pages 3–19, 2002. 1

[71] S. Romdhani and T. Vetter. Estimating 3D shape and texture using pixel intensity, edges, specular highlights, texture constraints and a prior. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2, pages 986–993, 2005. 2

[72] Zeyu Ruan, Changqing Zou, Longhai Wu, Gangshan Wu, and Limin Wang. SADRNet: Self-aligned dual face regression networks for robust 3D dense face alignment and reconstruction. *Transactions on Image Processing*, 30:5793–5806, 2021. 2

[73] James A Russell. A circumplex model of affect. *Journal of personality and social psychology*, 39(6):1161, 1980. 3

[74] Georgia Sandbach, Stefanos Zafeiriou, Maja Pantic, and Lijun Yin. Static and dynamic 3D facial expression recognition: A comprehensive survey. *Image and Vision Computing*, 30(10):683–697, 2012. 3

[75] Soubhik Sanyal, Timo Bolkart, Haiwen Feng, and Michael Black. Learning to regress 3D face shape and expression from an image without 3D supervision. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7763–7772, 2019. 1, 2, 5

[76] Matan Sela, Elad Richardson, and Ron Kimmel. Unrestricted facial geometry reconstruction using image-to-image translation. In *International Conference on Computer Vision (ICCV)*, pages 1576–1585, 2017. 2

[77] Caifeng Shan, Shaogang Gong, and Peter W McOwan. Facial expression recognition based on local binary patterns: A comprehensive study. *Image and vision Computing*, 27(6):803–816, 2009. 3

[78] Jiaxiang Shang, Tianwei Shen, Shiwei Li, Lei Zhou, Mingmin Zhen, Tian Fang, and Long Quan. Self-supervised monocular 3D face reconstruction by occlusion-aware multi-view geometry consistency. In *European Conference on Computer Vision (ECCV)*, volume 12360, pages 53–70, 2020. 2, 6, 7, 16

11

[79] Yingyan Shi, Qiaosha Zou, and Yiyun Zhang. Pose-robust facial expression recognition by 3D morphable model learning. In *International Conference on Computer and Communications (ICCC)*, pages 2458–2462, 2020. 3

[80] Ian Sneddon, Margaret McRorie, Gary Mckeown, and Jennifer Hanratty. The belfast induced natural emotion database. *Transactions on Affective Computing*, 3:32–41, 08 2013. 2

[81] Attila Szabó, Givi Meishvili, and Paolo Favaro. Unsupervised generative 3D shape learning from natural images. *CoRR*, abs/1910.00287, 2019. 2

[82] Ayush Tewari, Florian Bernard, Pablo Garrido, Gaurav Bharaj, Mohamed Elgharib, Hans-Peter Seidel, Patrick Pérez, Michael Zollhöfer, and Christian Theobalt. FML: face model learning from videos. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10812–10822, 2019. 1, 2

[83] Ayush Tewari, Mohamed Elgharib, Gaurav Bharaj, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zollhöfer, and Christian Theobalt. StyleRig: Rigging stylegan for 3D control over portrait images. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6141–6150, 2020. 1

[84] Ayush Tewari, Michael Zollhöfer, Pablo Garrido, Florian Bernard, Hyeongwoo Kim, Patrick Pérez, and Christian Theobalt. Self-supervised multi-level face model learning for monocular reconstruction at over 250 Hz. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2549–2559, 2018. 2

[85] Ayush Tewari, Michael Zollhöfer, Hyeongwoo Kim, Pablo Garrido, Florian Bernard, Patrick Perez, and Christian Theobalt. MoFA: model-based deep convolutional face autoencoder for unsupervised monocular reconstruction. In *International Conference on Computer Vision (ICCV)*, pages 1274–1283, 2017. 2

[86] Justus Thies, Michael Zollhöfer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Face2Face: Real-time face capture and reenactment of RGB videos. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2387–2395, 2016. 1, 2

[87] Ying-li Tian, Takeo Kanade, and Jeffrey F. Cohn. Recognizing action units for facial expression analysis. *Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 23(2):97–115, 2001. 3

[88] Antoine Toisoul, Jean Kossaifi, Adrian Bulat, Georgios Tzimiropoulos, and Maja Pantic. Estimation of continuous valence and arousal levels from faces in naturalistic conditions. *Nature Machine Intelligence*, 3(1):42–50, 2021. 2, 4, 5, 6, 7, 8, 13, 14, 21, 25

[89] Anh Tuan Tran, Tal Hassner, Iacopo Masi, and Gerard Medioni. Regressing robust and discriminative 3D morphable models with a very deep neural network. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1599–1608, 2017. 2

[90] Anh Tuan Tran, Tal Hassner, Iacopo Masi, Eran Paz, Yuval Nirkin, and Gérard Medioni. Extreme 3D face reconstruction: Seeing through occlusions. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3935–3944, 2018. 2

[91] Luan Tran, Feng Liu, and Xiaoming Liu. Towards high-fidelity nonlinear 3D face morphable model. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1126–1135, 2019. 2

[92] Thomas Vetter and Volker Blanz. Estimating coloured 3D face models from single images: An example based approach. In *European Conference on Computer Vision (ECCV)*, pages 499–513, 1998. 2

[93] Kaisiyuan Wang, Qianyi Wu, Linsen Song, Zhuoqian Yang, Wayne Wu, Chen Qian, Ran He, Yu Qiao, and Chen Change Loy. MEAD: A large-scale audio-visual dataset for emotional talking-face generation. In *European Conference on Computer Vision (ECCV)*, pages 700–717, 2020. 2

[94] Mei Wang, Weihong Deng, Jiani Hu, Xunqiang Tao, and Yaohai Huang. Racial faces in the wild: Reducing racial bias by information maximization adaptation network. In *International Conference on Computer Vision (ICCV)*, October 2019. 2

[95] Huawei Wei, Shuang Liang, and Yichen Wei. 3D dense face alignment via graph convolution networks. *arXiv preprint arXiv:1904.05562*, 2019. 2

[96] Zhen Wen and Thomas S. Huang. Capturing subtle facial motions in 3D face tracking. In *International Conference on Computer Vision (ICCV)*, pages 1343–1350, 2003. 3

[97] Shangzhe Wu, Christian Rupprecht, and Andrea Vedaldi. Unsupervised learning of probably symmetric deformable 3D objects from images in the wild. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–10, 2020. 2

[98] Haotian Yang, Hao Zhu, Yanru Wang, Mingkai Huang, Qiu Shen, Ruigang Yang, and Xun Cao. FaceScape: a large-scale high quality 3D face dataset and detailed riggable 3D face prediction. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 601–610, 2020. 2

[99] Xiaoxing Zeng, Xiaojiang Peng, and Yu Qiao. DF2Net: A dense-fine-finer network for detailed 3D face reconstruction. In *International Conference on Computer Vision (ICCV)*, 2019. 2

[100] Xiangyu Zhu, Zhen Lei, Xiaoming Liu, Hailin Shi, and Stan Z. Li. Face alignment across large poses: A 3D solution. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 146–155, 2016. 2

[101] Michael Zollhöfer, Justus Thies, Pablo Garrido, Derek Bradley, Thabo Beeler, Patrick Pérez, Marc Stamminger, Matthias Nießner, and Christian Theobalt. State of the art on monocular 3D face reconstruction, tracking, and applications. *Computer Graphics Forum*, 2018. 2

## A. Appendix

**Discussion of novelty:** In Sup. Mat. we aim to shed more light on the process that eventually led to EMOCA and the challenges that had to be overcome. The idea of using deep perceptual losses to supervise face reconstruction is not new. A critic might argue, that the novelty of EMOCA is very limited for exactly that reason. However, the fact remains that previous SOTA methods have a clear limitation when it comes to reconstructing faces that communicate the correct emotional content. And from the knowledge of this limitation, we conceived the idea of leveraging emotion recognition, an idea not previously attempted by any work on face reconstruction. The inventive novelty in EMOCA was in coming up with the idea in the first place. This idea, once explained, makes such an intuitive sense, it may lead the reader into thinking it is a straightforward change to an already functioning system. The idea, although very simple and elegant, was by no means easy to get to work and this is what we aim to explain next.

**Designing EMOCA:** Our work starts with the simple idea - how can we employ the findings from emotion recognition to improve face reconstruction? Leveraging a pretrained SOTA network for emotion recognition, similarly to the way face recognition networks were used seems like a natural choice. However, using its final outputs such as the expression class and valence and arousal levels is not sufficient. Clearly, these very low-dimensional labels, while they do carry some information about the emotional content, they likely exhibit a lot of ambiguity and are not sufficient to supervise 3D shapes. For instance, an expression classified as happy can take on many different shapes (a subtle smile, a big smile with an open mouth, an "inverted" smile, etc.) and similar reasoning could be applied for any other expression and for any levels of valence and arousal as well. Hence, these labels most likely do not provide a sufficient supervision signal for geometry. The next logical design choice is to leverage high dimensional deep features from a pretrained emotion recognition network. This choice can only make sense if the emotion feature in question is a "well-behaved" embedding space. Ideally we want similar features to represent faces of similar expressions and vice versa. Therefore, we conducted an emotion retrieval experiment, using a pretrained publicly available EmoNet model [88] and nearest neighbors search. This experiment is discussed in Sec. H. Having verified, that similar emotion features retrieve images of geometrically and semantically similar expressions, the next thing to be verified is whether the emotion feature carries a signal that is strong enough, to be utilizable for 3D reconstruction. This was particularly challenging and we comment on this further in Sec. D. Finally, having demonstrated that the emotion recognition features indeed carry enough information in order to supervise the geometry, we can finally incorporate the emotion consistency loss into a face reconstruction framework, arriving at EMOCA. In addition to the ablations listed in the main paper, we also add ablations on different architectures and weights for the emotion consistency loss in Sec. F

## B. Implementation details

**Emotion recognition metrics:** In the main paper, we evaluate emotion metrics in the same setting as Toisoul et al. [88]. The metrics are defines as follows RMSE stands for root mean squared error:

$$\text{RMSE}(Y, \hat{Y}) = \sqrt{\mathbb{E}[(Y - \hat{Y})^2]}.$$

SAGR stands for sign agreement and it evaluates whether the predicted value has the same sign as the ground truth:

$$\text{SAGR}(Y, \hat{Y}) = \frac{1}{n} \sum_{i=1}^{n} \delta\left(\text{sign}\left(y_i\right), \text{sign}\left(\hat{y}_i\right)\right).$$

Pearson correlation coefficient (PCC) measures the correlation between predictions and GT:

$$\text{PCC}(Y, \hat{Y}) = \frac{\mathbb{E}[(Y - \mu_Y)(\hat{Y} - \mu_{\hat{Y}})]}{\sigma_Y \sigma_{\hat{Y}}}.$$

Concordance correlation coefficient (CCC) incorporates the PCC but also penalizes signals which are still correlated according to PCC but have different means:

$$\text{CCC}(Y, \hat{Y}) = \frac{2\sigma_Y \sigma_{\hat{Y}} \text{PCC}(Y, \hat{Y})}{\sigma_Y^2 + \sigma_{\hat{Y}}^2 + \left(\mu_Y - \mu_{\hat{Y}}\right)^2}.$$

**Emotion recognition loss function:** We train our emotion networks with the same loss function as defined by Toisoul et al. [88].

$$\mathcal{L}_{\text{categories}}(Y, \hat{Y}) = \text{Cross entropy}(Y, \hat{Y}) = -\sum_{i=1}^{n} \hat{y}_i \log\left(y_i\right)$$

The complete loss function for emotion recognition is then defined as:

$$\mathcal{L}(Y, \hat{Y}) = \mathcal{L}_{\text{categories}}(Y, \hat{Y}) + \frac{\alpha}{\alpha + \beta + \gamma} \mathcal{L}_{\text{MSE}}(Y, \hat{Y})$$
$$+ \frac{\beta}{\alpha + \beta + \gamma} \mathcal{L}_{\text{PCC}}(Y, \hat{Y}) + \frac{\gamma}{\alpha + \beta + \gamma} \mathcal{L}_{\text{CCC}}(Y, \hat{Y}),$$

where $\alpha$, $\beta$ and $\gamma$ are shake-shake regularization coefficients [31] uniformly sampled from the interval $[0, 1]$ for each training batch and:

$$\mathcal{L}_{\text{MSE}}(Y, \hat{Y}) = \text{MSE}_{\text{valence}}(Y, \hat{Y}) + \text{MSE}_{\text{arousal}}(Y, \hat{Y})$$

$$\mathcal{L}_{\text{PCC}}(Y, \hat{Y}) = 1 - \frac{\text{PCC}_{\text{valence}}(Y, \hat{Y}) + \text{PCC}_{\text{arousal}}(Y, \hat{Y})}{2}$$

$$\mathcal{L}_{\text{CCC}}(Y, \hat{Y}) = 1 - \frac{\text{CCC}_{\text{valence}}(Y, \hat{Y}) + \text{CCC}_{\text{arousal}}(Y, \hat{Y})}{2}.$$

Unlike the work of Toisoul et al. [88], we do not use knowledge distillation as its improvements are marginal and make the training process much more complex.

**Image-based emotion recognition:** We investigate emotion recognition networks based on different architectures, ResNet-50 [37], Swin Transformer [55], and EmoNet [88]. We train all models on AffectNet [59], using the training/validation/test split proposed by Toisoul et al. [88]. The ResNet-50 and Swin Transformer based models are pretrained on ImageNet [18]. During training, the training images are sampled such that each of the 7 expression labels appears with the same frequency. This sampling is crucial to maximize the performance of the emotion networks, as the AffectNet training set is not balanced. We use the Adam optimizer with learning rate of 0.0001, $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The batch size used for training is 64. Each model is trained for a maximum of 20 epochs with early stopping, and the model with the lowest validation error is selected.

**3DMM-based emotion recognition:**

In Section 5.2 of the paper (Tab. 1) and Table 3, we evaluate different face reconstruction methods by recognizing emotions from the regressed 3DMM parameters. Specifically, we train a 4-layer MLP with Batch Normalization and LeakyReLUs to output valence and arousal levels and expression classes from the regressed identity and expression parameters (see Figs. 5 and 6 for details). The size of each hidden layer is 2048. We train the 3DMM-based recognition on AffectNet similarly to the image-based emotion recognition. The loss function is identical to the one used for image-based emotion recognition. The batch size used for training is 64. We use the Adam optimizer with a learning rate of 0.0001, $\beta_1 = 0.9$ and $\beta_2 = 0.999$.

**Detail stage training:** The detail stage training follows the training protocol of DECA [28]. The coarse model part is kept fixed, while detail encoder and decoder are trained. This stage uses VGGFace2 [14] and VoxCeleb2 [17] images, due to the necessity of having multiple images per identity. We optimize following losses: photometric loss, ID-MRF perceptual loss which encourages reconstruction of higher frequency detail (compared to the coarse mesh), as well as the soft symmetry loss and displacement regularization. Further, to disentangle identity and expression dependent details, we employ DECA's detail consistency loss, where each batch contains k images of each subject, and the detail codes are exchanged randomly between the predictions for each identity For our training, we set k=3 and batch size of 4 identities, totalling 12 input images per batch. For more details, see the original DECA publication.

## C. Qualitative evaluation

In addition to the performance in emotion analysis on the AffectNet dataset in the main paper, we also test EMOCA on AFEW-VA [49]. The results are reported in Tab. 3.

## D. Emotion optimization

We can use our emotion consistency loss for additional tasks. Here we consider the problem of expression retargeting. Given two face images, a source identity image $I_S$ and a target expression image $I_T$ of potentially two different people with different expressions, poses, cameras, and lighting, our goal is to optimize for the (unknown) target expression $\hat{\psi}_T$. Formally, we infer the FLAME parameters $E_c(I_S)$ and $E_c(I_T)$ for both images. Then, with some abuse of notation, we render $I_R(\psi) = R(M(\beta_S, \theta_T, \psi), \alpha_S, \mathbf{l}_T, \mathbf{c}_T)$, the FLAME mesh with source identity shape $\beta_S$, source albedo $\alpha_S$, and target pose $\theta_T$, target camera $\mathbf{c}_T$, target lighting $\mathbf{l}_T$, and the optimization expression parameters $\psi$. We then extract the emotion features of the rendering $\epsilon_R(\psi) = A(I_R(\psi))$ and the target image $\epsilon_T = A(I_T)$, and optimize:

$$\hat{\psi}_T = \arg\min_{\psi} d(\epsilon_R(\psi), \epsilon_T) + \lambda_\psi L_\psi, \qquad (9)$$

with $d(\epsilon_1, \epsilon_2) = \|\epsilon_1 - \epsilon_2\|_2$, expression regularizer $L_\psi = \|\psi\|_2^2$, and regularizer weight $\lambda_\psi = 1e{-}3$. We use gradient descent for the optimization. Below we show optimization results, and an analysis of the convergence and sensitivity to the initialization.

**On Emotion Network Architecture:** Figure 7 shows emotion optimization results using different emotion recognition network. This indicates that the original released EmoNet is not suitable for emotion optimization. Instead, we use the ResNet-50 architecture as default model.

**On Initialization:** Figure 8 further shows the influence of the initialization on the optimized emotion. These results demonstrate that 3DMMs, when rendered, can in fact be animated with a deep perceptual emotion similarity loss.

**On Jaw Optimization:** A perceptive reader may ask, why we optimize only for the expression parameters $\psi$ and not also for the jaw pose $\theta_{jaw}$. After all, the jaw position most certainly has an effect on the perceived emotion. We have struggled with the jaw optimization issue for quite a long time, unable to get acceptable results as the jaw pose parameter optimization makes this optimization unstable - the jaw would always be posed to an unrealistic or at least very incorrect pose. Fixing the jaw pose to a reasonable estimate however (such as DECA's prediction) makes the optimization stable and produces good results. We hypothesise that this instability could be caused by the following:

1. FLAME is missing a comprehensive prior for the jaw pose. We experimented with simplistic hand-crafted
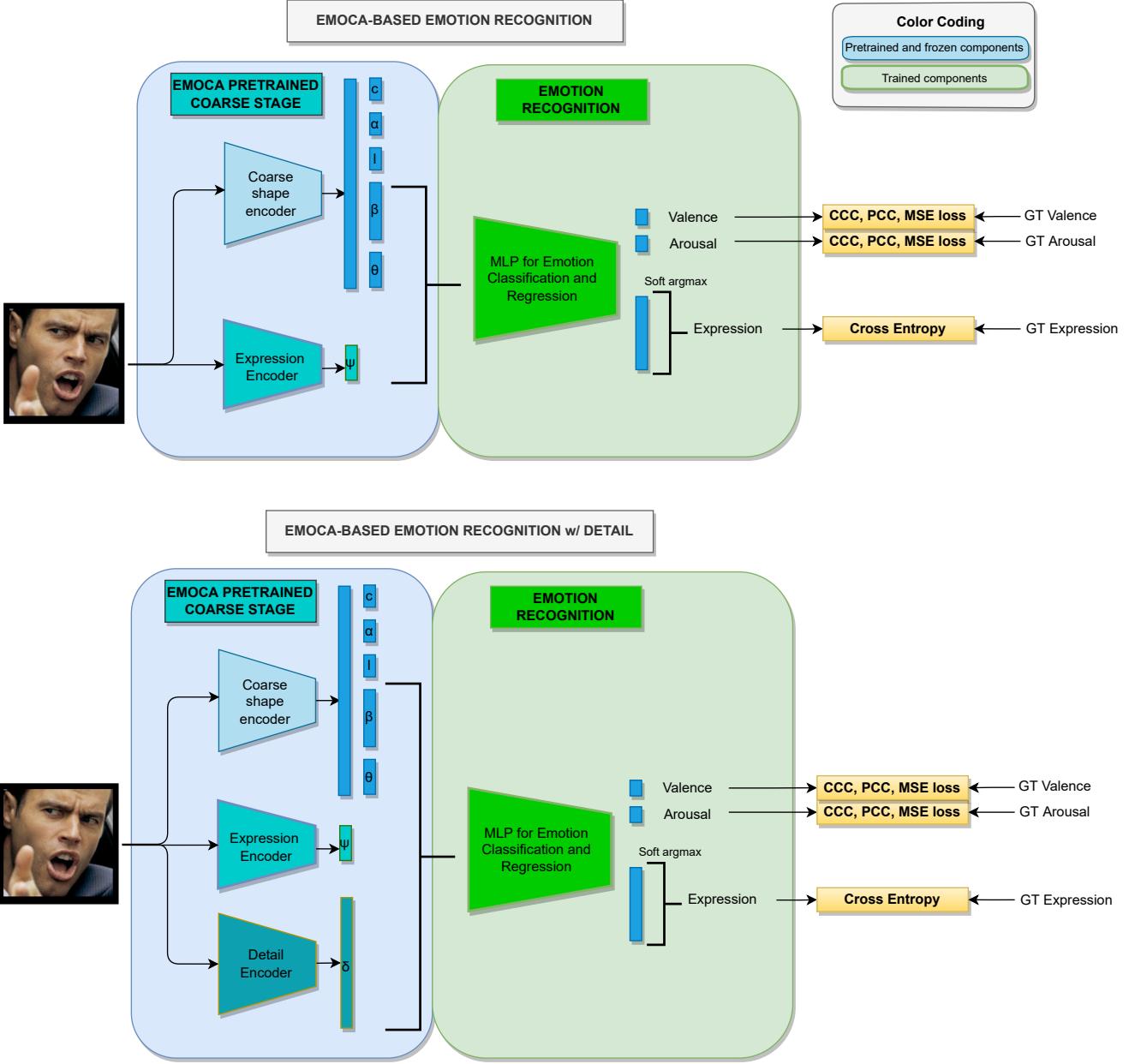
Figure 5. The architecture of EMOCA-based emotion recognition. Top: EMOCA emotion recognition with coarse parameters. From the pretrained coarse stage we extract the shape parameters $\boldsymbol{\beta}$, expression parameters $\boldsymbol{\psi}$ and jaw pose $\boldsymbol{\theta}_{jaw}$. A similar approach is taken for DECA-based recognition, except that DECA does not have a dedicated expression encoder. These are fed to an MLP to regress valence and arousal and classify expression. Bottom: emotion recognition for EMOCA-based reconstruction methods with detail code included.

priors (such as distance or squared distance from the expected pose) but this did not yield any improvement. It is possible that the creating a more comprehensive prior (other than the Gaussian prior for FLAME's expression space), a prior that entangles the expression and jaw pose spaces is necessary. This makes for an interesting direction for future work.

2. Emotion optimization involves optimizing a deep feature vector and while we have demonstrated that similar emotion features belong to similar expressions, we have not eliminated the possibility, that the emotion network can be "attacked" to produce the desired features with a distorted images. An optimization process, in which the jaw is not fixed could results in an adversarial attack on the network that forces it to produce a
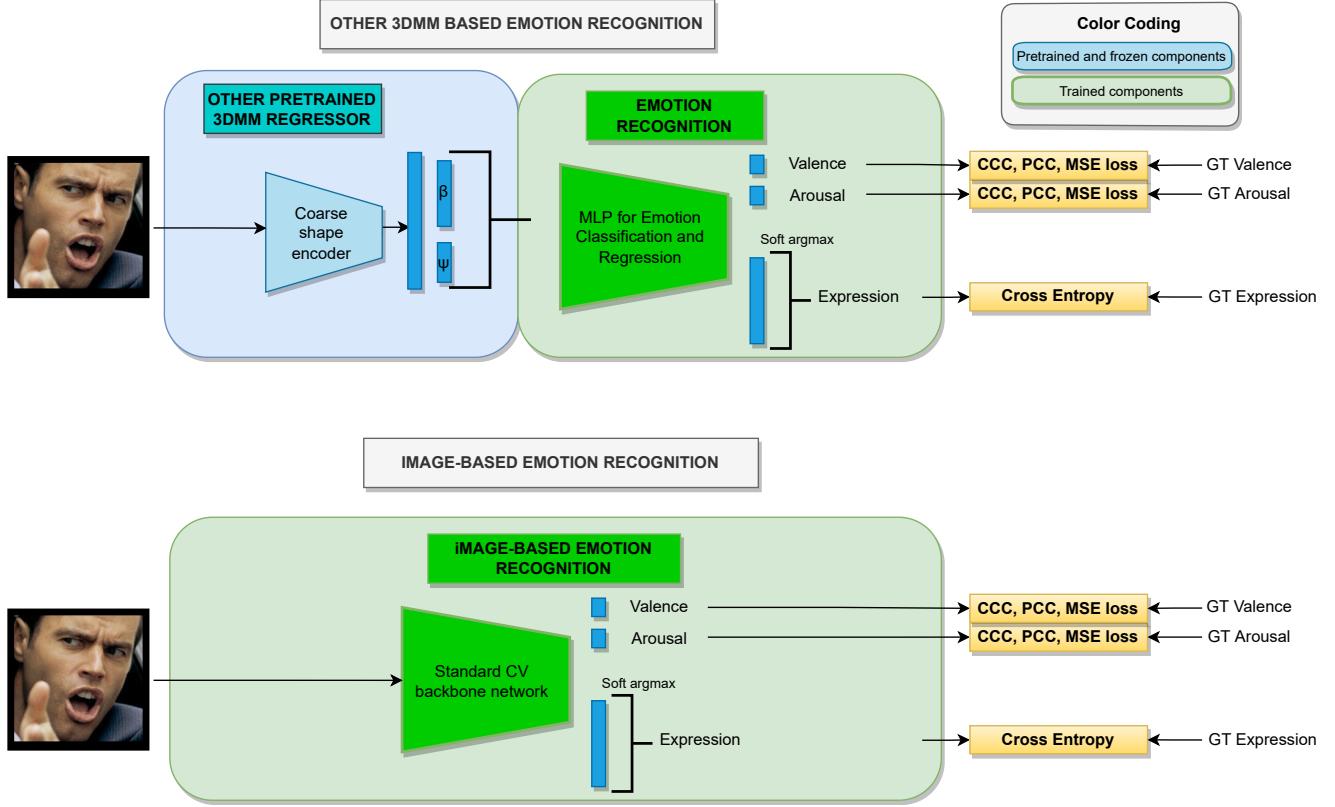
**Figure 6.** The architecture of other emotion recognition netowrks. Top: emotion recognition for other 3DMM-based reconstruction methods (Deep3DFace [20], 3DDFA-V2 [36], MGCNet [78]. These have a single decoder that regress to the Basel Face Model [64] parameter space, which does not model jaw pose explicitly. Therefore only $\beta$ and $\psi$ are considered. Bottom: a standard image-based network trained for emotion recognition. Both types of emotion recognition are trained with the same supervision.

| Model | V-PCC ↑ | V-CCC ↑ | V-RMSE ↓ | V-SAGR ↑ | A-PCC ↑ | A-CCC ↑ | A-RMSE ↓ | A-SAGR ↑ |
|---|---|---|---|---|---|---|---|---|
| EmoNet | 0.59 | 0.54 | 0.22 | 0.61 | 0.55 | 0.49 | *0.22* | *0.80* |
| Deep3DFace | *0.64* | *0.59* | *0.21* | **0.65** | 0.55 | 0.48 | **0.21** | **0.81** |
| ExpNet | 0.31 | 0.25 | 0.27 | 0.55 | 0.36 | 0.30 | 0.24 | 0.79 |
| MGCNet | 0.54 | 0.50 | 0.23 | 0.62 | 0.49 | 0.44 | 0.23 | 0.79 |
| 3DDFA | 0.41 | 0.38 | 0.27 | 0.57 | 0.44 | 0.41 | 0.24 | 0.78 |
| DECA (coarse) | 0.57 | 0.53 | 0.23 | 0.62 | *0.55* | *0.50* | 0.22 | **0.81** |
| DECA /w detail | 0.57 | 0.53 | 0.23 | 0.63 | 0.53 | 0.49 | *0.22* | *0.80* |
| EMOCA (Ours) | **0.65** | **0.63** | *0.21* | 0.64 | **0.57** | **0.54** | 0.22 | *0.80* |
| EMOCA /w detail (Ours) | **0.68** | **0.65** | **0.20** | 0.64 | *0.56* | *0.53* | 0.22 | *0.80* |

**Table 3.** Emotion recognition performance on AFEW-VA [49]. All emotion regressors are pretrained on AffectNet and finetuned on the AFEW-VA using 5-fold Cross-Validation (CV). The reported numbers are averaged across the 5-fold CV runs. EMOCA performs best, followed by Deep3DFace. Surprisingly, both of these methods outperform EmoNet. Other 3D-based methods follow.

similar emotion feature vector.

# E. Perceptual study

Section 5.2 of the paper evaluates the amount emotion conveyed by the reconstructed 3D geometry in a perceptual study. Figure 9 gives the full confusion matrix of the participants' labels of real images (rows) and the labels of the reconstructions (columns). Figure 10 further compares the ground truth emotion labels with the participants' classifications of the reconstructions. For completeness, we also include the confusion matrix of participants' labeling of the real images in Fig. 11.
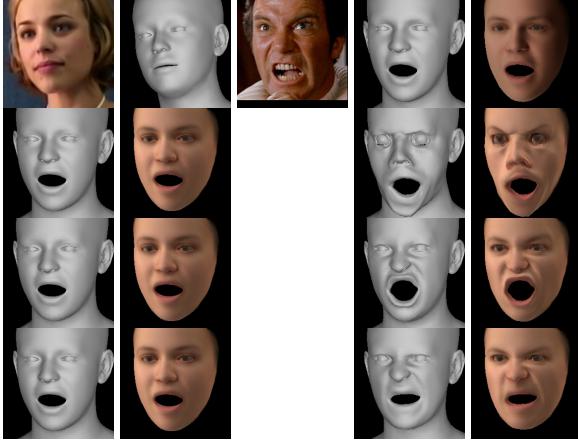
Figure 7. Emotion optimization example. The first row contains a source image, its DECA reconstruction, a target image, its DECA reconstruction, and the colored reconstruction. The following rows contain: the initialization of the optimization w/o and w/ color (left) and the optimization result w/o and w/ color (right). The different rows use different emotion recognition networks for optimization. The second row uses the original released EmoNet, the third row a self-trained EmoNet, and the bottom row using our ResNet-50 model. While EmoNet gives SOTA emotion recognition results, it is less suitable for our task of emotion-driven expression optimization or reconstruction.

## F. Emotion consistency

**Emotion network architecture:** The choice of architecture for emotion supervision plays a critical role. While all architectures perform comparatively well on the emotion recognition task, they are not equally suitable as supervision for our 3D face reconstruction task. Fig. 12 visually compares EMOCA models trained with different emotion recognition networks as supervision. Again, the SOTA emotion recognition architecture - EmoNet, is not suitable as it produces unacceptable artifacts. Furthermore, the SWIN [55] transformer backbone, which is considered to be superior to the ResNet [37] architecture, also produces some undesirable artifacts. Hence, the ResNet backbone was used for the final model of the emotion recognition network.

**Emotion consistency weight:** We have experimented with different values of the emotion consistency loss weight term $\lambda_{emo}$. This is a crucial factor of successfully training EMOCA. If the weight is too small, the emotion is not captured well enough. At the same time, high values lead to unnaturally over-exaggerated expressions. A visual ablation of this phenomenon can be found in Fig. 13 and Fig. 14 for two different emotion network architectures; ResNet-50 [37] and SWIN-B [55].

**Additional ablations:** We further evaluate the impact of the similarity metric used for the emotion similarity, the effect



Figure 8. Sensitivity of the emotion optimization to initialization. The first row contains a source image, its DECA reconstruction, a target image, its DECA reconstruction, and the colored reconstruction. The following rows contain: the initialization of the optimization w/o and w/ color (left) and the optimization result w/o and w/ color (right). Note that the optimization process is only modifying the expression coefficients $\psi$ and not the jaw rotation $\boldsymbol{\theta}_{jaw}$. While the process usually converges to meaningful results, the most favorable outcome is obtained, when initializing the process with the target expression coefficients $\psi$ and pose $\boldsymbol{\theta}$, which correspond to the second row.

of adding a landmark reprojection error to the loss function, and the effect of the relative landmark losses (mouth closure, lip corner distance and eye closure). Finally, we analyze the effect of using DECA's training data instead of AffectNet. You can see the results in Fig. 15.

## G. Emotional retargeting

EMOCA regresses FLAME [53] parameters and expression dependent geometric details. The disentanglement of the coarse identity and expression geometry and the identity and expression dependent details allows us to animate EMOCA's reconstructions. We demonstrate this by animat-
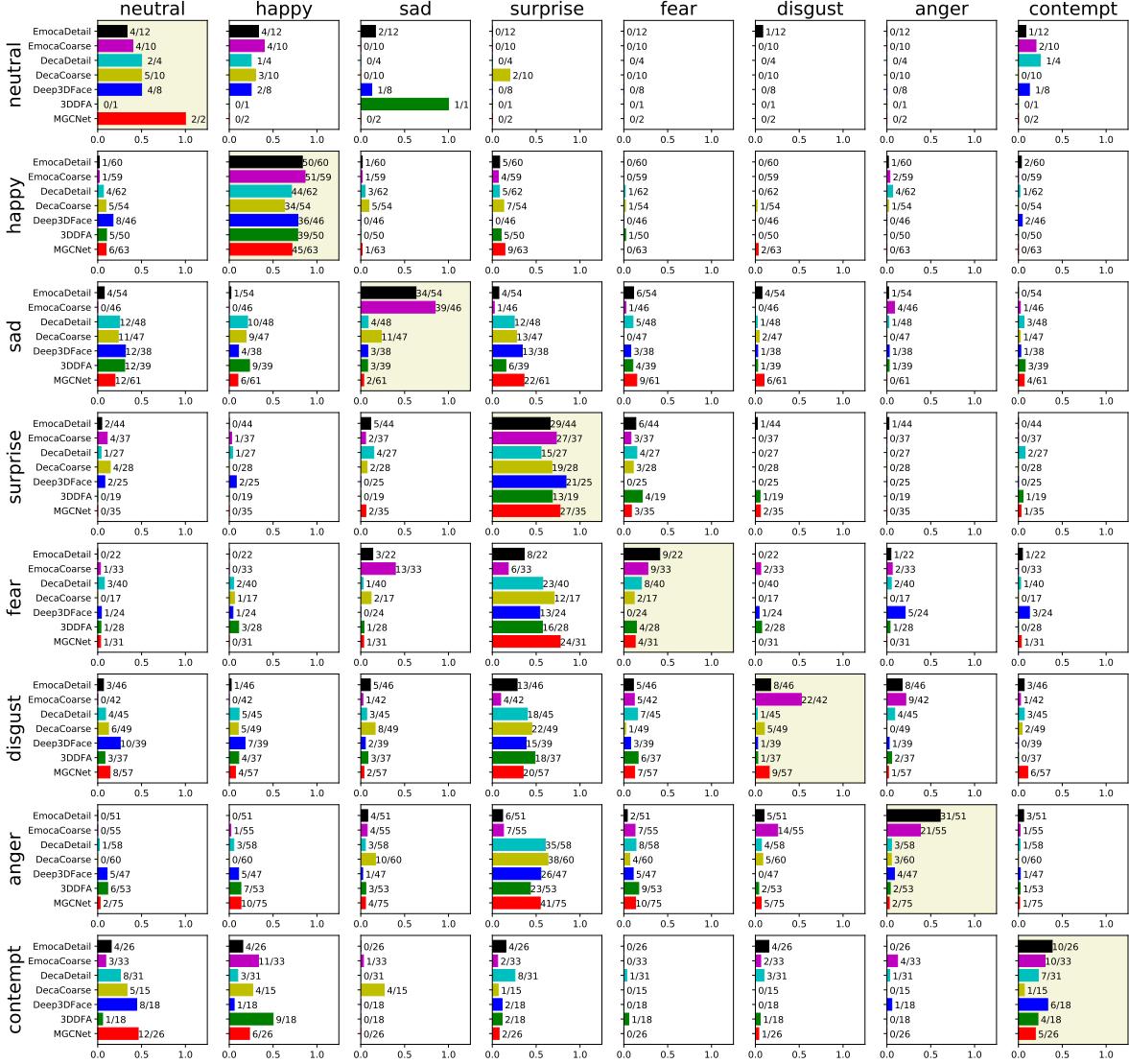
Figure 9. This figure contains the confusion matrices of participant's labels of the real image and the reconstructed images for each method. The x-axis of each cell gives the ratio of participants' reconstruction labels and real image labels and the absolute number is written next to each bar. The accuracy of each method for a particular expression class is on the diagonal. You can see that both variants of EMOCA (detail and coarse) are superior to the other methods. Furthermore, off-diagonal you can observe how the label of meshes reconstructed by EMOCA is much less confused for other labels, compared to other methods. Finally, the confusion matrix highlights how other methods are not capable of producing expressions of fear, disgust and anger. Instead these are confused with surprise. EMOCA does not suffer from the same limitation. However, participants did have some trouble distinguishing reconstructions of disgust and anger. Please note that the first row (neutral) shows a small number of samples. This happens because our perceptual study did not contain neutral images.
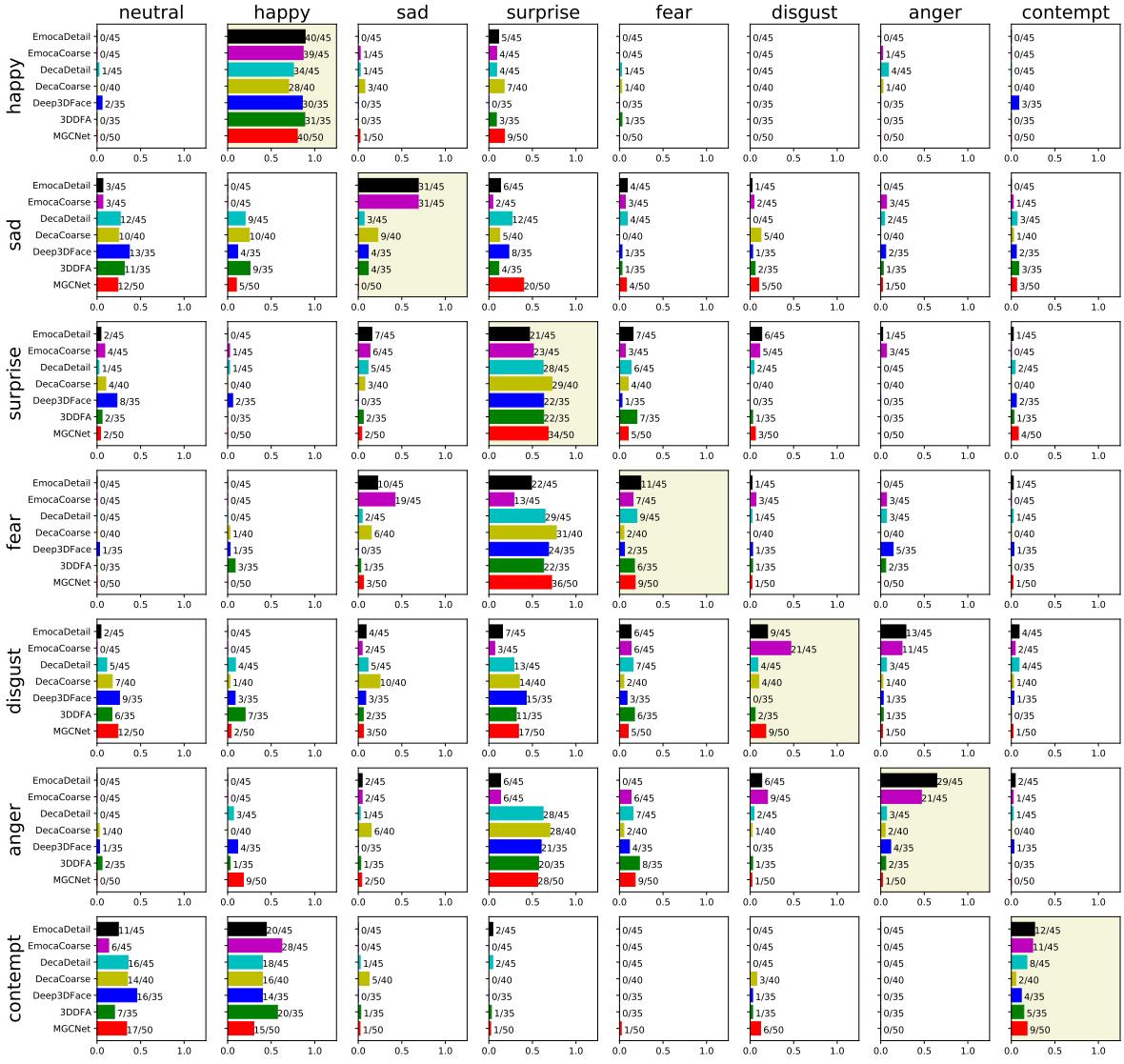
18

Figure 10. This figure contains the confusion matrices of participant's labels of the reconstructions w.r.t. to the ground truth labels (as opposed to users' subjective labels, which you can find in Fig. 9. Please note that neutral expressions were not given in the study, which is why the matrix only has six rows (neutral excluded).

ing a source 3D face using a video sequence of another actor. Figure 16 demonstrates two things, first, EMOCA reconstructions convey emotions of the source images, and second, the animated faces of other subjects convey the

same emotion. The emotional fidelity hence is preserved in the animated face of the other subject.

Figure 11. This figure contains the confusion matrices of participant's labels of the real images w.r.t. to ground truth images. While this figure does not compare the performance of methods, it serves as a baseline comparison to Fig. 10. Classifying expression is subjective. While our participants mostly agreed with our ground truth, there were disagreements for the negatively charged expressions of fear, disgust, anger and particularly contempt.

## H. Emotion retrieval

Our work relies on the following key hypothesis. The emotion recognition networks learn a useful embedding of emotion. The following properties are desirable:

- Images of faces with similar expressions conveying similar emotions are close in this embedding space.
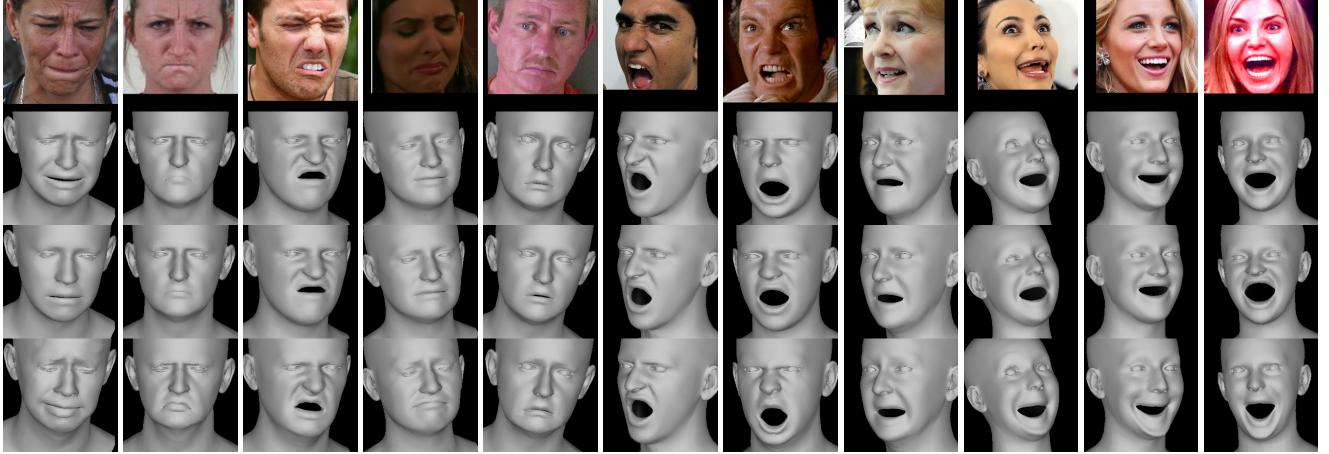
Figure 12. Comparison of different EMOCA models, supervised by different emotion networks. From top to bottom: ResNet-50 [37], SWIN-B [55], EmoNet [88]. All three networks affect the reconstruction in different ways. EMOCA-ResNet produces the best visual results and is our model of choice. EMOCA-SWIN produces results of slightly lower visual quality. Finally, EMOCA-EmoNet sometimes produces unrealistic expressions, which makes EmoNet less suitable for this task.



Figure 13. Comparison of models trained with different weights of the emotion consistency loss $\lambda_{emo}$. The emotion network used was ResNet-50 [37]. Top row consists of input images. Different values of $\lambda_{emo}$ follow. From top to bottom 0, 0.1, 0.5, 1 (final EMOCA), 5, 10.

- Images of faces with dissimilar expressions/emotions are farther apart in this space.

- Invariance to pose, identity and lighting and back-ground.

We employ the publicly released model of EmoNet [88] and use the 256-dimensional feature output of the last con-

Figure 14. Comparison of models trained with different weights of the emotion consistency loss $\lambda_{emo}$. The emotion network used was SWIN-B [55]. Top row consists of input images. Different values of $\lambda_{emo}$ follow. From top to bottom 0, 0.1, 0.5, 1, 5, 10. While SWIN-B suffers from fewer artifacts compared to ResNet-50 when changing the weight, we have deemed the visual quality of results produce by a ResNet-supervised EMOCA slightly better, which is why ResNet was selected for the final model.

volutional layer as emotion embedding. We then extract the emotion embedding for faces in the Aff-Wild2 video dataset [48]. For the emotion retrieval given an image, we seek the nearest neighbors w.r.t. L2 distance metric in the dataset. Figure 17 shows the 10 nearest neighbors for multiple images. For comparison, we repeat the process for the ground truth (GT) valence and arousal labels of the Aff-Wild2 dataset in Fig. 18.
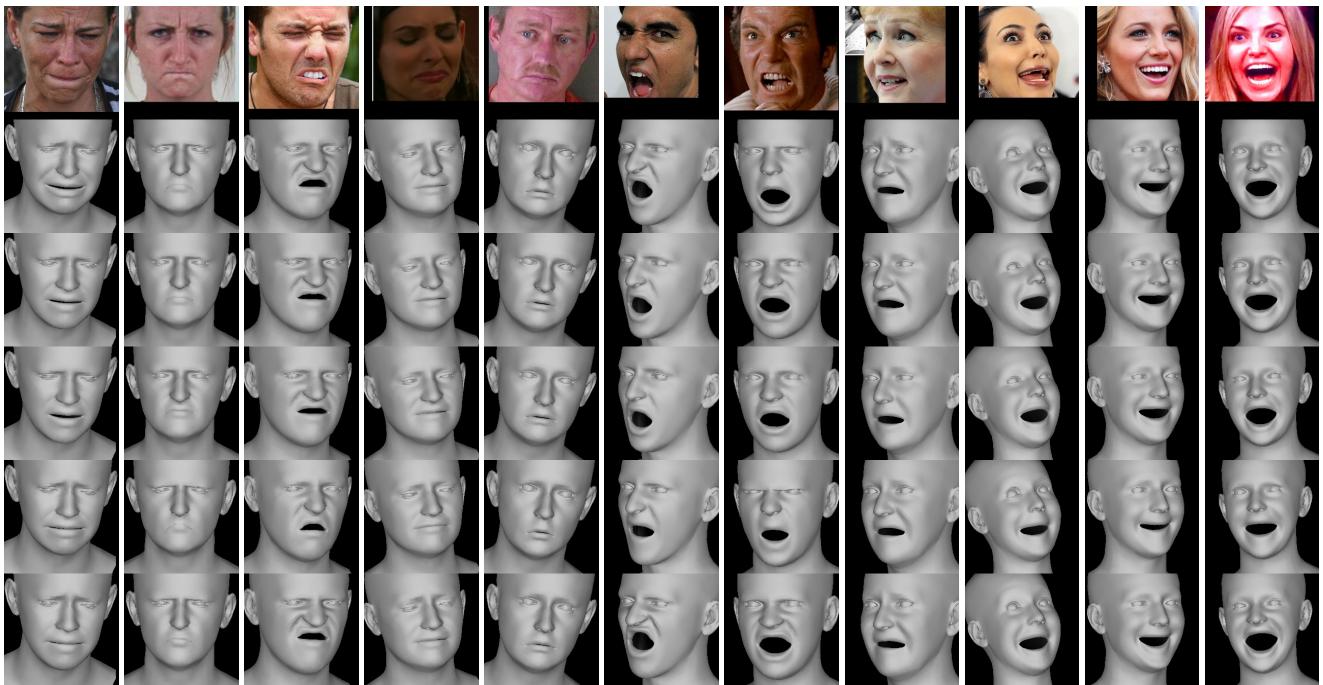
Figure 15. A visual comparison of model with different changes. First row consists of input images. The next three rows use different metrics for evaluating emotion similarity - L2 (EMOCA), L1 and cosine similarity. As you can observe, the selection of the metric is not critical for performance. The following row drops the relative landmark losses (mouth closure, eye closure and lip corner distance). Observe that this has a negative effect on the samples, particularly the mouth region. Final row is EMOCA model trained on the same data as DECA instead of AffectNet. You can see that it achieves a very similar result compared to EMOCA trained on AffectNet. This highlights an interesting finding - once an emotion recognition network has been trained, it can be used for supervision even on datasets that do not strictly guarantee a balanced representation of emotional states, such as face recognition datasets.
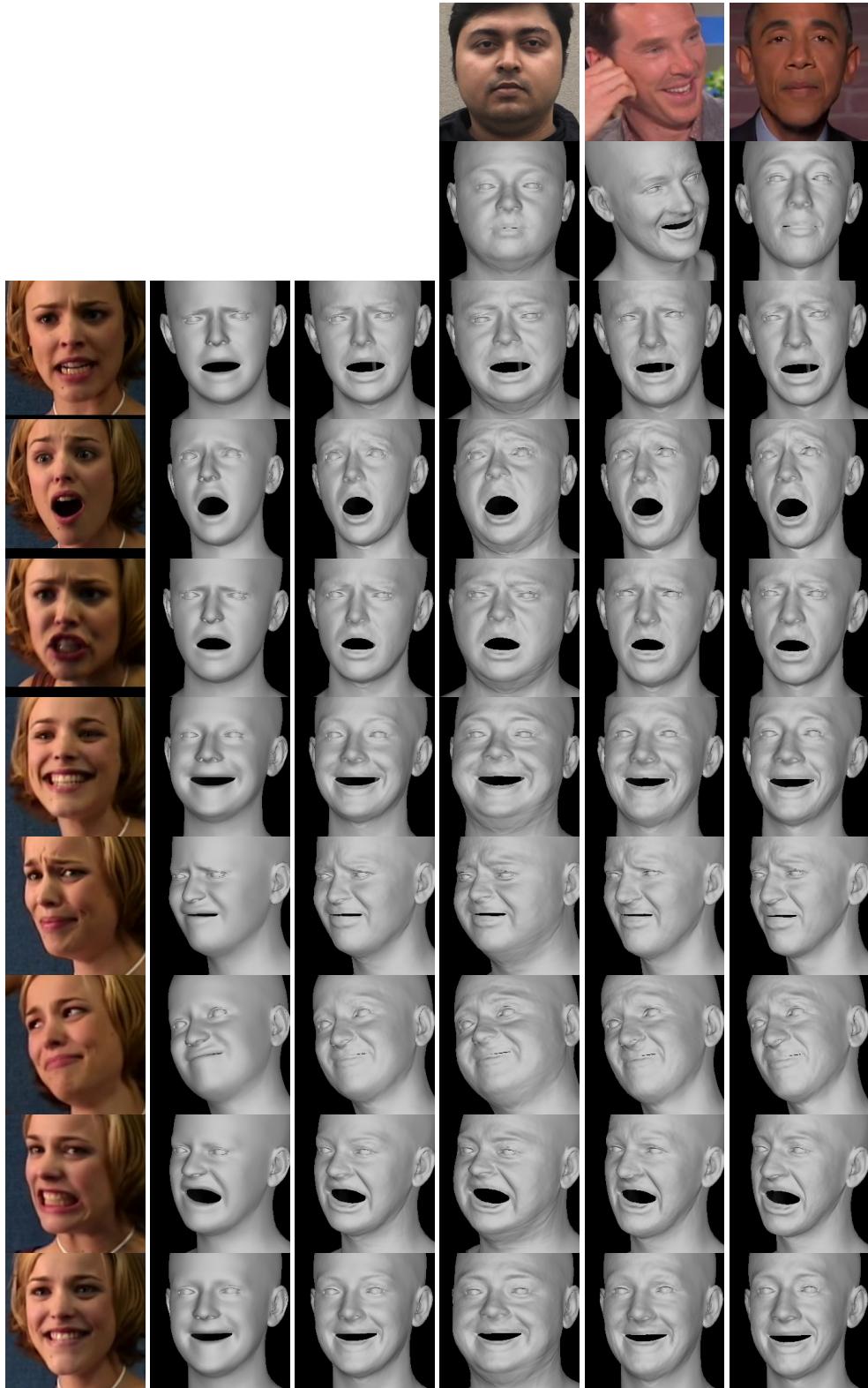
Figure 16. **Emotional retargetting.** From left to right. The input image, coarse reconstruction, detailed reconstruction, emotion retargeted to the coarse identity above. Observe that while the identity and the person-specific detailed displacements change with the source actor, the emotion fidelity is preserved. For the entire sequence in motion, please see the supplementary video.

Figure 17. Examples of nearest neighbor retrieval using the EmoNet [88] feature. We searched for up to 100 neighbors. We only include up to 1 NN per video to avoid retrieving consecutive frames. Left: query image, Right: ordered nearest neighbors from different clips. Observe how all of the retrieved faces communicate very similar emotional content.
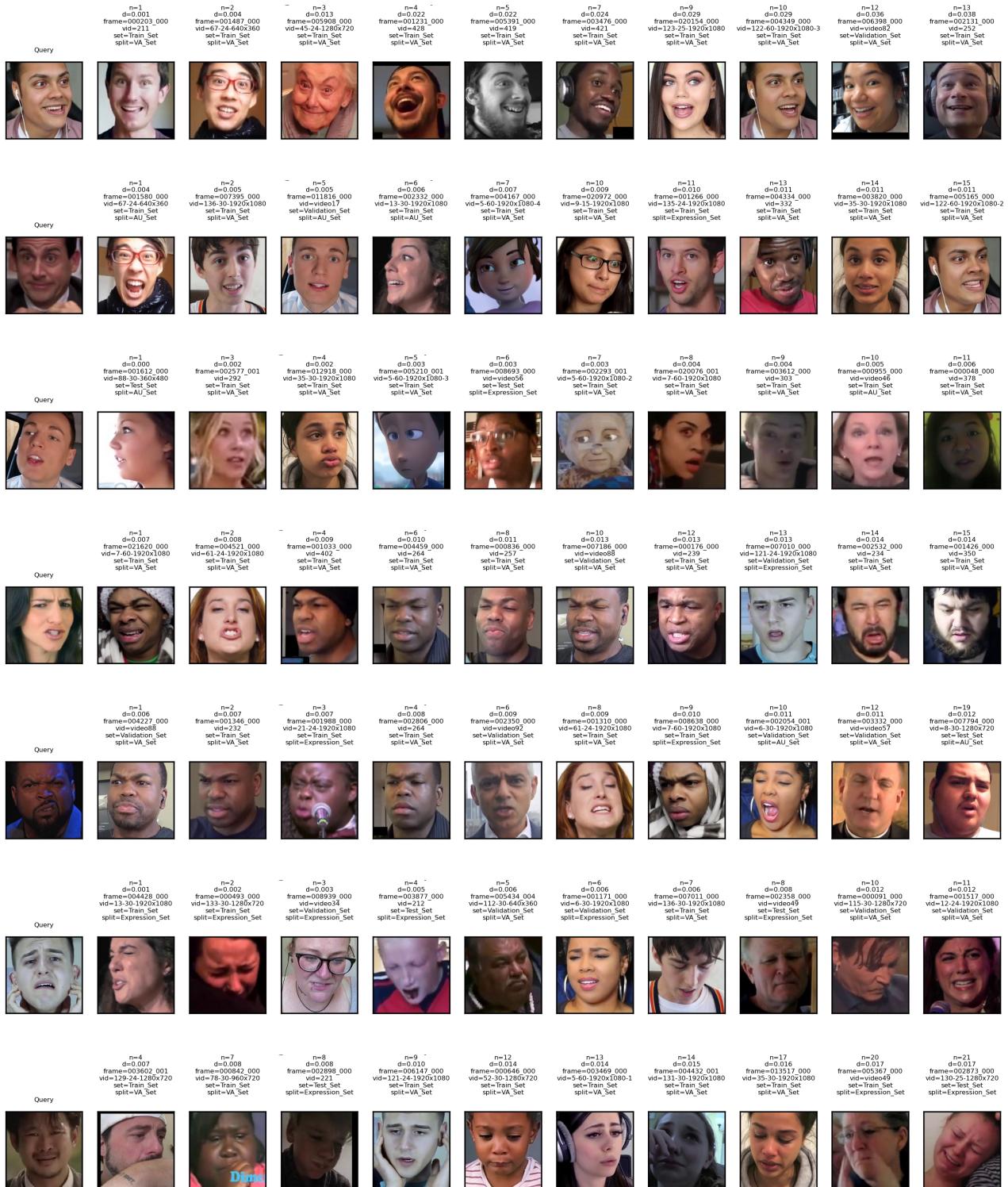
Figure 18. Examples of nearest neighbor retrieval using the ground truth annotated valence and arousal space on the AffWild2 [48] dataset. While the retrieved faces do have some degree of similarity, the quality of retrieval compared to the EmoNet feature is lower.