# Sparse to Dense Dynamic 3D Facial Expression Generation

Naima Otberdout

Univ. Lille, CNRS, Inria, Centrale Lille, UMR 9189 CRIStAL, F-59000 Lille, France

`naima.otberdout@univ-lille.fr`

Claudio Ferrari

Deparment of Architecture and Engineering

University of Parma, Italy

`claudio.ferrari2@unipr.it`

Mohamed Daoudi

IMT Nord Europe, Institut Mines-Télécom, Univ. Lille, Centre for Digital Systems, F-59000 Lille, France
Univ. Lille, CNRS, Centrale Lille, Institut Mines-Télécom, UMR 9189 CRIStAL, F-59000 Lille, France

`mohamed.daoudi@imt-nord-europe.fr`

Stefano Berretti
Media Integration ad Communication Center
University of Florence, Italy

`stefano.berretti@unifi.it`

Alberto Del Bimbo
Media Integration ad Communication Center
University of Florence, Italy

`alberto.delbimbo@unifi.it`

## Abstract

*In this paper, we propose a solution to the task of generating dynamic 3D facial expressions from a neutral 3D face and an expression label. This involves solving two sub-problems: (i) modeling the temporal dynamics of expressions, and (ii) deforming the neutral mesh to obtain the expressive counterpart. We represent the temporal evolution of expressions using the motion of a sparse set of 3D landmarks that we learn to generate by training a manifold-valued GAN (Motion3DGAN). To better encode the expression-induced deformation and disentangle it from the identity information, the generated motion is represented as per-frame displacement from a neutral configuration. To generate the expressive meshes, we train a Sparse2Dense mesh Decoder (S2D-Dec) that maps the landmark displacements to a dense, per-vertex displacement. This allows us to learn how the motion of a sparse set of landmarks influences the deformation of the overall face surface, independently from the identity. Experimental results on the CoMA and D3DFACS datasets show that our solution brings significant improvements with respect to previous solutions in terms of both dynamic expression generation and mesh reconstruction, while retaining good generalization to unseen data. The code and the pretrained*

*model will be made publicly available.*

## 1. Introduction

Synthesizing dynamic 3D (4D) facial expressions aims at generating realistic face instances with varying expressions or speech-related movements that dynamically evolve across time, starting from a face in neutral expression. It finds application in a wide range of graphics applications spanning from 3D face modeling, to augmented and virtual reality for animated films and computer games. While recent advances in generative neural networks have made possible the development of effective solutions that operate on 2D images [17, 37], the literature on the problem of generating facial animation in 3D is still quite limited.

To perform a faithful and accurate 3D facial animation, three main challenges arise. First, the identity of the subject whose neutral face is used as starting point for the sequence should be maintained across time. Second, the applied deformation should correspond to the specified expression/motion that is provided as input, and should be applicable to any neutral 3D face. Incidentally, these are major challenges in 3D face modeling, which require disentangling structural face elements related to the identity, *e.g.*, nose or jaw shape, from deformations related to the
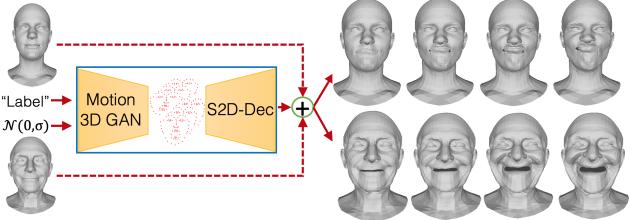
Figure 1. **3D dynamic facial expression generation**: A GAN generates the motion of 3D landmarks from an expression label and noise; A decoder expands the animation from the landmarks to a dense mesh, while keeping the identity of a neutral 3D face,

movable face parts, *e.g.*, mouth opening/closing. Finally, it is required to model the temporal dynamics of the specified expression so to obtain realistic animations.

Some previous works tackled the problem by capturing the facial expression of a subject frame-by-frame and transferring it to a target model [7]. However, in this case the temporal evolution is not explicitly modeled, so the problem reduces to transferring a tracked expression to a neutral 3D face. Some other works animated a 3D face mesh given an arbitrary speech signal and a static 3D face mesh as input [13, 29]. Also in this case, the temporal evolution is guided by an external input, similar to a tracked expression. Instead, here we are interested in animating a face just starting from a neutral face and an expression label.

In our solution, which is illustrated in Figure 1, the temporal evolution and the mesh deformation are decoupled and modeled separately in two network architectures. A manifold-valued GAN (*Motion3DGAN*) accounts for the expression dynamics by generating a temporally consistent motion of 3D landmarks corresponding to the input label from noise. The landmarks motion is encoded using the *Square Root Velocity Function* (SRVF) and compactly represented as a point on a hypersphere. Then, a Sparse2Dense mesh Decoder (*S2D-Dec*) generates a dense 3D face guided by the landmarks motion for each frame of the sequence. To effectively disentangle identity and expression components, the landmarks motion is represented as a per-frame displacement from a neutral configuration. Instead of directly generating a mesh, the S2D-Dec expands the landmarks displacement to a dense, per-vertex displacement, which is finally used to deform the neutral mesh. The intuition that led to this architecture is the following: the movement induced on the face surface by the underlying facial muscles is consistent across subjects. In addition, it causes the vertex motion to be locally correlated as muscles are smooth surfaces. We thus train the decoder to learn how the displacement of a sparse set of points influences the displacement of the whole face surface. This has the advantage that structural face parts, *e.g.*, nose or forehead, which are not influenced by facial expressions are not deformed, helping

in maintaining the identity traits stable. Furthermore, the network can focus on learning expressions at a fine-grained level of detail and generalize to unseen identities.

In summary, the main contributions of our work are: *(i)* we propose an original method to generate dynamic sequences of 3D expressive scans given a neutral 3D mesh and an expression label. Our approach has the capability of generating strong and diverse expression sequences, with high generalization ability to unseen identities and expressions; *(ii)* we adapt a specific GAN architecture [37] for dynamic 3D landmarks generation, and design a decoder for expressive mesh reconstruction from a neutral mesh and landmarks. Differently from common auto-encoders, the proposed S2D-Dec learns to generate a per-vertex displacement map from a few control points, allowing accurate mesh deformations where the structural face parts remain stable; *(iii)* we designed a novel reconstruction loss that weighs the contribution of each vertex based on its distance from the landmarks. This proved to augment the decoder capability of generating accurate expressions.

## 2. Related Work

Our work is related to methods for 3D face modeling, facial expression generation guided by landmarks, and dynamic generation of 3D faces, *i.e.*, 4D face generation.

**3D face modeling**. The 3D Morphable face Model (3DMM) as originally proposed in [2] is the most popular solution for modeling 3D faces. The original model and its variants [3, 6, 19, 21, 34, 36, 38] capture face shape variations both for identity and expression based on linear formulations, thus incurring in limited modeling capabilities. For this reason, non-linear encoder-decoder architectures are attracting more and more attention. This comes at the cost of reformulating convolution and pooling/unpooling like operations on the non-regular mesh support [5, 33, 42]. Ranjan *et al.* [40] proposed an auto-encoder architecture that builds upon newly defined spectral convolution operators, and pooling operations to down-/up-sample the mesh. Bouritsas *et al.* [4] improved upon the above by proposing a novel graph convolutional operator enforcing consistent local orderings on the vertices of the graph through the *spiral operator* [32]. Despite their impressive modeling precision, a recent work [19] showed that they heavily suffer from poor generalization to unseen identities. This limits their practical use in tasks such as face fitting or expression transfer. We finally mention that other approaches exist to learn generative 3D face models, such as [1, 35]. However, instead of dealing with meshes they use alternative representations for 3D data, such as depth images or UV-maps.

To overcome the above limitation, we go beyond self-reconstruction and propose a mesh decoder that, differently from previous models, learns expression-specific mesh deformations from a sparse set of landmark displacements.

**Facial expression generation guided by landmarks**. Recent advances in neural networks made facial landmark detection reliable and accurate both in 2D [10, 15, 43] and 3D [24, 46]. Landmarks and their motion are a viable way to account for facial deformations as they reduce the complexity of the visual data, and have been commonly used in several 3D face related tasks, *e.g.*, reconstruction [21, 31] or reenactment [18, 22]. Despite some effort was put in developing landmark-free solutions for 3D face modeling [8, 9, 23], some recent works investigated their use to model the dynamics of expressions. Wang *et al.* [44] proposed a framework that decouples facial expression dynamics, encoded into landmarks, and face appearance using a conditional recurrent network. Otberdout *et al.* [37] proposed an approach for generating videos of the six basic expressions given a neutral face image. The geometry is captured by modeling the motion of landmarks with a GAN that learns the distribution of expression dynamics.

These methods demonstrated the potential of using landmarks to model the dynamics of expressions and generate 2D videos. In our work, we instead tackle the problem of modeling the dynamics in 3D, exploring the use of the motion of 3D landmarks to both model the temporal evolution of expressions and animate a 3D face.

**4D face generation**. While many researchers tackled the problem of 3D mesh deformation, the task of 3D facial motion synthesis is yet more challenging. A few studies addressed this issue by exploiting audio features [29, 45], speech signal [13] or tracked facial expressions [7] to generate facial motions. However, none of these explicitly model the temporal dynamics and resort to external information.

To the best of our knowledge, the work in [39] is the only approach that specifically addressed the problem of dynamic 3D expression generation. In that framework, the motion dynamics is modeled with a temporal encoder based on an LSTM, which produces a per-frame latent code starting from a per-frame expression label. The codes are then fed to a mesh decoder that, similarly to our approach, generates a per-vertex displacement that is summed to a neutral 3D face to obtain the expressive meshes. Despite the promising results reported in [39], we identified some limitations. First, the LSTM is deterministic and for a given label the exact same displacements are generated. Our solution instead achieves diversity in the output sequences by generating from noise. Moreover, in [39] the mesh decoder generates the displacements from the latent codes, making it dependent from the temporal encoder. In our solution, the motion dynamics and mesh displacement generation are decoupled, using landmarks to link the two modules. The S2D-Dec is thus independent from Motion3DGAN, and can be used to generate static meshes as well given a arbitrary set of 3D landmarks as input. This permits us to use the decoder for other tasks such as expression/speech transfer.

Finally, as pointed out in [39], the model cannot perform extreme variations well. Using landmarks allowed us to define a novel reconstruction loss that weighs the error of each vertex with respect to its distance from the landmarks, encouraging accurate modeling of the movable parts. Thanks to this, we are capable of accurately reproducing from slight to strong expressions, and generalize to unseen motions.

## 3. Proposed Method

Our approach consists of two specialized networks as summarized in Figure 2. Motion3DGAN accounts for the temporal dynamics and generates the motion of a sparse set of 3D landmarks from noise, provided an expression label, *e.g.*, happy, angry. The motion is represented as per frame landmark displacements with respect to a neutral configuration. These displacements are fed to a decoder network, S2D-Dec, that constructs the dense point-cloud displacements from the sparse displacements given by the landmarks. These dense displacements are then added to a neutral 3D face to generate a sequence of expressive 3D faces corresponding to the initial expression label. In the following, we separately describe the two networks.

### 3.1. Generating Sparse Dynamic 3D Expressions

Facial landmarks were shown to well encode the temporal evolution of facial expressions [28, 37]. Motivated by this fact, we generate the facial expression dynamics based on the motion of 3D facial landmarks. Given a set of $k$ 3D landmarks, $Z(t) = (x_i(t), y_i(t), z_i(t))_{i=1}^{k}$, with $Z(0)$ being the neutral configuration, their motion can be seen as a trajectory in $\mathbb{R}^{k \times 3}$ and can be formulated as a parameterized curve in $\mathbb{R}^{k \times 3}$ space. Let $\alpha : I = [0, 1] \to \mathbb{R}^{k \times 3}$ represent the parameterized curve, where each $\alpha(t) \in \mathbb{R}^{k \times 3}$. For the purpose of modeling and studying our curves, we adopt the Square-Root Velocity Function (SRVF) proposed in [41]. The SRVF $q(t) : I \to \mathbb{R}^{k \times 3}$ is defined by:

$$q(t) = \frac{\dot{\alpha}(t)}{\sqrt{\|\dot{\alpha}(t)\|_2}}, \tag{1}$$

with the convention that $q(t) = 0$ if $\dot{\alpha}(t) = 0$. This function proved effective for tasks such as human action recognition [14] or 3D face recognition [16]. Similar to this work, Otberdout *et al.* [37] proposed to use the SRVF representation to model the temporal evolution of 2D facial landmarks, which makes it possible to learn the distribution of these points and generate new 2D facial expression motions. In this paper, we extend this idea to 3D by proposing the Motion3DGAN model to generate the motion of 3D facial landmarks represented using the SRVF encoding in (1).

Following [37], we remove the scale variability of the resulting motions by scaling the $\mathbb{L}^2$-norm of these functions to 1. As a result, we transform the motion of 3D
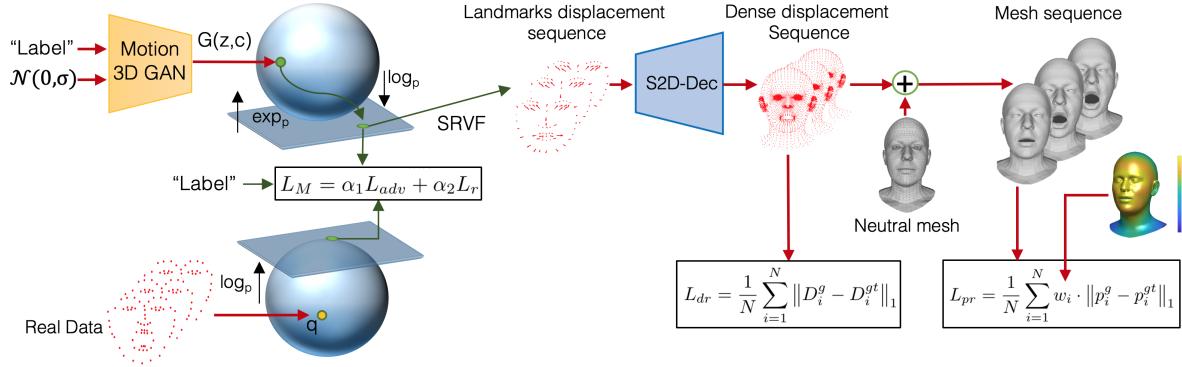
Figure 2. **Overview of our framework.** Motion3DGAN generates the motion $q(t)$ of 3D landmarks corresponding to an expression label from a noise vector $z$. The module is trained guided by a reconstruction loss $L_r$ and adversarial loss $L_{adv}$. The motion $q(t)$ is converted to a sequence of landmark displacements $d_i$, which are fed to S2D-Dec. From each $d_i$, the decoder generates a dense displacement $D_i^g$. A neutral mesh is then summed to the dense displacements to generate the expressive meshes $\mathbf{S}^g$. S2D-Dec is trained under the guidance of a displacement loss $L_{dr}$ and our proposed weighted reconstruction loss $L_{pr}$.

facial landmarks to points on a hypersphere $\mathcal{C} = \{q : [0,1] \to \mathbb{R}^{k \times 3}, \|q\| = 1\}$. The resulting representations are manifold-valued data that cannot be handled with traditional generative models.

To learn the distribution of the SRVF representations, we propose Motion3DGAN as an extension of Motion-GAN [37], a conditional version of the Wasserstein GAN for manifold-valued data [27]. It maps a random vector $z$ to a point on the hypersphere $\mathcal{C}$ conditioned on an input class label. Motion3DGAN is composed of two networks trained adversarially: a generator $G$ that learns the distribution of the 3D landmark motions, and a discriminator $D$ that distinguishes between real and generated 3D landmark motions. Motion3DGAN is trained by a weighted sum of an adversarial loss $L_{adv}$ and a reconstruction loss $L_r$ such that $L_M = \alpha_1 L_{adv} + \alpha_2 L_r$. The former is given by:

$$
\begin{aligned}
L_{adv} = \quad & \mathbb{E}_{q \sim \mathbb{P}_q} \left[ D\left(\log_p(q), c\right) \right] \\
& - \mathbb{E}_{z \sim \mathbb{P}_z} \left[ D\left(\log_p\left(\exp_p(G(z,c))\right)\right) \right] \\
& + \lambda E_{\hat{q} \sim \mathbb{P}_{\hat{q}}} \left[ \left( \|\nabla_{\hat{q}} D(\hat{q})\|_2 - 1 \right)^2 \right].
\end{aligned} \quad (2)
$$

In (2), $q \sim \mathbb{P}_q$ is an SRVF sample from the training set, $c$ is the expression label (*e.g.*, mouth open, eyebrow) that we encode as a one-hot vector and concatenate to a random noise $z \sim \mathbb{P}_z$. The last term of the adversarial loss represents the gradient penalty of the Wasserstein GAN [25]. Specifically, $\hat{q} \sim \mathbb{P}_{\hat{q}}$ is a random point sampled uniformly along straight lines between pairs of points sampled from $\mathbb{P}_q$ and the generated distribution $\mathbb{P}_g$:

$$
\hat{q} = (1 - \tau) \log_p(q) + \tau \log_p(\exp_p(G(z,c))), \quad (3)
$$

where $0 \le \tau \le 1$, and $\nabla_{\hat{q}} D(\hat{q})$ is the gradient *w.r.t.* $\hat{q}$. The functions $\log_p(.)$ and $\exp_p(.)$ are the logarithm and the exponential maps, respectively, defined in a particular point

$p$ of the hypersphere. They map the SRVF data forth and back to a tangent space of $\mathcal{C}$ (details in the supplementary material). Finally, the reconstruction loss is defined as:

$$
L_r = \|\log_p(\exp_p(G(z,c))) - \log_p(q)\|_1, \quad (4)
$$

where $\|.\|_1$, represents the $L_1$-norm, and $q$ is the ground truth SRVF corresponding to the condition $c$. The generator and discriminator architectures are similar to [37].

The SRVF representation is reversible, which makes it possible to recover the curve $\alpha(t)$ from a new generated SRVF $q(t)$ by,

$$
\alpha(t) = \int_0^t \|q(s)\| q(s) ds + \alpha(0). \quad (5)
$$

Where $\alpha(0)$ represents the initial landmarks configuration $Z(0)$. Accordingly, using this equation, we can apply the generated motion to *any* landmark configuration, making it robust to identity changes.

### 3.2. From Sparse to Dense 3D Facial Expressions

Our final goal is to animate the neutral mesh $\mathbf{S}^n$ to obtain a novel 3D face $\mathbf{S}^g$ reproducing some expression, yet maintaining the identity structure of $\mathbf{S}^n$. Given this, we point at generating the displacements of the mesh vertices from the sparse displacements of the landmarks to animate $\mathbf{S}^n$. In the following, we assume all the meshes have a fixed topology, and are in full point-to-point correspondence.

Let $\mathcal{L} = \left\{ \left( \mathbf{S}_1^n, \mathbf{S}_1^{gt}, Z_1^n, Z_1^{gt} \right), \ldots, \left( \mathbf{S}_m^n, \mathbf{S}_m^{gt}, Z_m^n, Z_m^{gt} \right) \right\}$ be the training set, where $\mathbf{S}_i^n = (p_1^n, \ldots, p_N^n) \in \mathbb{R}^{N \times 3}$ is a neutral 3D face, $\mathbf{S}_i^{gt} = (p_1^{gt}, \ldots, p_N^{gt}) \in \mathbb{R}^{N \times 3}$ is a 3D expressive face, $Z_i^n \in \mathbb{R}^{k \times 3}$ and $Z_i^{gt} \in \mathbb{R}^{k \times 3}$ are the 3D landmarks corresponding to $\mathbf{S}_i^n$ and $\mathbf{S}_i^{gt}$, respectively. We transform this set to a training set of sparse and dense

displacements, $\mathcal{L}' = \{(D_1, d_1), \ldots, (D_m, d_m)\}$ such that, $D_i = \mathbf{S}_i^{gt} - \mathbf{S}_i^n$ and $d_i = Z_i^{gt} - Z_i^n$. Our goal here is to find a mapping $h : \mathbb{R}^{k \times 3} \to \mathbb{R}^{N \times 3}$ such that $\mathbf{D}_i \approx h(d_i)$. We designed the function $h$ as a decoder network (S2D-Dec), where the mapping is between a sparse displacement of a set of landmarks and the dense displacement of the entire mesh points. Finally, in order to obtain the expressive mesh, the dense displacement map is summed to a 3D face in neutral expression, *i.e.*, $\mathbf{S}_i^e = \mathbf{S}_i^n + \mathbf{D}_i$. The S2D-Dec network is based on the spiral operator proposed in [4]. Our architecture includes five spiral convolution layers, each one followed by an up-sampling layer. More details on the architecture can be found in the supplementary material.

In order to train this network, we propose to use two different losses, one acting directly on the displacements and the other controlling the generated mesh. The reconstruction loss of the dense displacements is given by,

$$L_{dr} = \frac{1}{N} \sum_{i=1}^{N} \left\| D_i^g - D_i^{gt} \right\|_1, \tag{6}$$

where $D^g$ and $D^{gt}$ are the generated and the ground truth dense displacements, respectively. To further improve the reconstruction accuracy, we add a loss that minimizes the error between $\mathbf{S}^g$ and the ground truth expressive mesh $\mathbf{S}^{gt}$. We observed that vertices close to the landmarks are subject to stronger deformations. Other regions like the forehead, instead, are relatively stable. To give more importance to those regions, we defined a weighted version of the $L1$ loss:

$$L_{pr} = \frac{1}{N} \sum_{i=1}^{N} w_i \cdot \left\| p_i^g - p_i^{gt} \right\|_1. \tag{7}$$

We define the weights as the inverse of the Euclidean distance of each vertex $p_i$ in the mesh from its closest landmark $Z_j$, *i.e.* $w_i = \frac{1}{\min d(p_i, Z_j)}$, $\forall j$. This provides a coarse indication of how much each $p_i$ contributes to the expression generation. Since the mesh topology is fixed, we can precompute the weights $w_i$ and re-use them for each sample. Weights are then re-scaled so that they lie in $[0, 1]$. Vertices corresponding to the landmarks, *i.e.*, $p_i = Z_j$ for some $j$, are hence assigned the maximum weight. We will show this strategy provides a significant improvement with respect to the standard $L1$ loss. The total loss used to train the S2D-Dec is given by $L_{S2D} = \beta_1.L_{dr} + \beta_2.L_{pr}$.

# 4. Experiments

We validated the proposed method in a broad set of experiments on two publicly available benchmark datasets.
**CoMA dataset** [40]: It is a common benchmark employed in other studies [4, 40]. It consists of 12 subjects, each one performing 12 extreme and asymmetric expressions. Each expression comes as a sequence of meshes $\mathbf{S} \in \mathbb{R}^{N \times 3}$ (140

meshes on average), with $N = 5,023$ vertices.
**D3DFACS dataset** [11]: We used the registered version of this dataset [30], which has the same topology of CoMA. It contains 10 subjects, each one performing a different number of facial expressions. In contrast to CoMA, this dataset is labeled with the activated action units of the performed facial expression. It is worthy to note that the expressions of D3DFACS are highly different from those in CoMA.

## 4.1. Training Details

In order to keep Motion3DGAN and S2D-Dec decoupled, they are trained separately. We used CoMA to train Motion3DGAN, since this dataset is labeled with facial expression classes. We manually divided each sequence into sub-sequences of length 30 starting from the neutral to the apex frame. The first sub-sequence for each subject and expression is used as test set. We use all the others for training as we generate from a random noise at test time.[1] Then, we encoded the motion of $k = 68$ landmarks from the sub-sequences in the SRVF representation, and used them to train Motion3DGAN. The landmarks were first centered and normalized to unit norm. Each of the 12 expression labels was encoded as a one-hot vector, concatenated with a random noise vector of size 128.

To comprehensively evaluate the capability of S2D-Dec of generalizing to either unseen identities or expressions, we performed subject-independent and expression-independent cross-validation experiments. For the subject-independent experiment, we used a 4-fold cross-validation protocol for CoMA, training on 9 and testing on 3 identities in each fold. On D3DFACS, we used the last 7 identities for training and the remaining 3 as test set. Concerning the expression-independent splitting, we used a 4-fold cross-validation protocol for CoMA, training on 9 and testing on 3 expressions in each fold. For D3DFACS, given the different number of expression per subject, the first 11 expressions were used for testing and trained on the rest.

We trained both Motion3DGAN and S2D-Dec using the Adam optimizer, with learning rate of 0.0001 and 0.001 and mini-batches of size 128 and 16, respectively. Motion3DGAN was trained for 8000 epochs, while 300 epochs were adopted for S2D-Dec. The hyper-parameters of the Motion3DGAN and S2D-Dec losses were set empirically to $\alpha_1 = 1$, $\alpha_2 = 10$, $\beta_1 = 1$ and $\beta_2 = 0.1$. We chose the mean SRVF of the CoMA data as a reference point $p$, where we defined the tangent space of $\mathcal{C}$.

## 4.2. 3D Expression Generation

For evaluation, we set up a baseline by first comparing against standard 3DMM-based fitting methods. Similar to

---

[1]For reproducibility, the list of the sub-sequences used to train Motion3DGAN will be publicly released.

previous works [20, 31], we fit $\mathbf{S}^n$ to the set of target landmarks $Z^e$ using the 3DMM components. Since the deformation is guided by the landmarks, we first need to select a corresponding set from $\mathbf{S}^n$ to be matched with $Z^e$. Given the fixed topology of the 3D faces, we can retrieve the landmark coordinates by indexing into the mesh, *i.e.*, $Z^n = \mathbf{S}^n(\mathbf{I}_z)$, where $\mathbf{I}_z \in \mathbb{N}^n$ are the indices of the vertices that correspond to the landmarks. We then find the optimal deformation coefficients that minimize the Euclidean error between the target landmarks $Z^e$ and the neutral ones $Z^n$, and use the coefficients to deform $\mathbf{S}^n$. In the literature, several 3DMM variants have been proposed. We experimented the standard PCA-based 3DMM and the DL-3DMM in [20]. We chose this latter variant as it is conceptually similar to our proposal, being constructed by learning a dictionary of deformation displacements. For fair comparison, we built the two 3DMMs using a number of deformation components comparable to the size of the S2D-Dec input, *i.e.*, $68 \times 3 = 204$. For PCA, we used either 38 components (retaining the $99\%$ of the variance) and 220, while for DL-3DMM we used 220 dictionary atoms.

With the goal of comparing against other deep models, we also considered the Neural3DMM [4]. It is a mesh autoencoder tailored for learning a non-linear latent space of face variations and reconstructing the input 3D faces. In order to compare it with our model, we modified the architecture and trained the model to generate an expressive mesh $\mathbf{S}^g$ given its neutral counterpart as input. To do so, we concatenated the landmarks displacement (of size 204) to the latent vector (of size 16) and trained the network towards minimizing the same $L_{pr}$ loss used in our model. All the compared methods were trained on the same data. Finally, we also identified the FLAME model [31]. Unfortunately, the training code is not available, and using the model pretrained on external data would not be a fair comparison.

The mean per-vertex Euclidean error between the generated meshes and their ground truth is used as standard performance measure, as in the majority of works [4,19,39,40]. Note that we exclude the Motion3DGAN model here as we do not have the corresponding ground-truth for the generated landmarks (they are generated from noise). Instead, we make use of the ground truth motion of the landmarks.

### 4.2.1 Comparison with Other Approaches

Table 1 shows a clear superiority of S2D-Dec over state-of-the-art methods for both the protocols and datasets, proving its ability to generate accurate expressive meshes close to the ground truth in both the case of unseen identities or expressions. In Figure 3, the cumulative per-vertex error distribution on the expression-independent splitting further highlights the precision of our approach, which can reconstruct 90%-98% of the vertices with an error lower than

| Method | Expression Split | | Identity Split | |
| --- | --- | --- | --- | --- |
| | CoMA | D3DFACS | CoMA | D3DFACS |
| PCA-220 | $0.76 \pm 0.73$ | $0.42 \pm 0.44$ | $0.80 \pm 0.73$ | $0.56 \pm 0.56$ |
| PCA-38 | $0.90 \pm 0.84$ | $0.44 \pm 0.45$ | $0.93 \pm 0.82$ | $0.58 \pm 0.56$ |
| DL3DMM [20] | $0,86 \pm 0,80$ | $0.73 \pm 1.15$ | $0.89 \pm 0.79$ | $1.15 \pm 1.50$ |
| Neural [4] | $0.75 \pm 0.85$ | $0.59 \pm 0.86$ | $3.74 \pm 2.34$ | $2.09 \pm 1.37$ |
| **Ours** | $\mathbf{0.52 \pm 0.59}$ | $\mathbf{0.28 \pm 0.31}$ | $\mathbf{0.55 \pm 0.62}$ | $\mathbf{0.27 \pm 0.30}$ |

Table 1. Reconstruction error (mm) on expression-independent (left) and identity-independent (right) splits: comparison with PCA-$k$ 3DMM ($k$ components), DL-3DMM (220 dictionary atoms), and Neural3DMM.
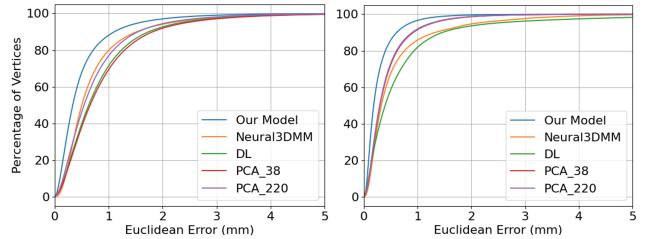


Figure 3. Cumulative per-vertex Euclidean error between PCA-based 3DMM models, DL-3DMM, Neural3DMM, and our proposed model, using expression-independent cross-validation on the CoMA (left) and D3DFACS (right) datasets.

$1mm$. While other fitting-based methods retain satisfactory precision in both the protocols, we note that the performance of Neural3DMM [4] significantly drop when unseen identities are considered. This outcome is consistent to that reported in [19], in which the low generalization ability of these models is highlighted. We also note that results for the identity-independent protocol were never reported in the original papers [4, 40]. Overall, our solution embraces the advantages of both approaches, being as general as fitting solutions yet more accurate.

Figure 4 shows some qualitative examples by reporting error heatmaps in comparison with PCA, DL-3DMM [20] and Neural3DMM [4] for the identity-independent splitting. The ability of our model as well as PCA and DL-3DMM to preserve the identity of the ground truth comes out clearly, in accordance with the results in Table 1. By contrast, Neural3DMM shows high error even for the neutral faces, which proves its inability to keep the identity of an unseen face. Indeed, differently from to the other methods, Neural3DMM encodes the neutral face in a latent space and predicts the 3D coordinates of the points directly, which introduces some changes on the identity of the input face. This evidences the efficacy of our S2D-Dec, that instead learns per-point displacements instead of point coordinates.
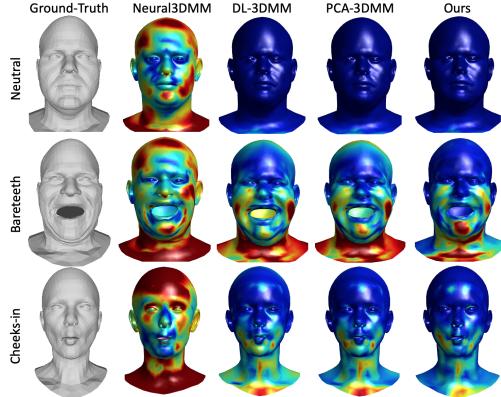
Figure 4. Mesh reconstruction error (red=high, blue=low) of our model and other methods.

| Method | Error (mm) |
|---|---|
| $S_1 : L_{dr}$ | $1.27 \pm 1.88$ |
| $S_2 : S_1 + L_{pr}$ w/o distance weights | $0.92 \pm 1.33$ |
| $S_3 : S_1 + L_{pr}$ | $0.50 \pm 0.56$ |

Table 2. Ablation study on the reconstruction loss of S2D-Dec.

#### 4.2.2 Ablation Study

We report here an ablation study to highlight the contribution of each loss used to train S2D-Dec, with particular focus on our proposed weighted-$L1$ reconstruction loss. We conducted this study on the CoMA dataset using the first three identities as a testing set and training on the rest. This evaluation is based on the mean per-vertex error between the generated and the ground truth meshes. We evaluated three baselines, $S1$, $S2$ and $S3$. For the first baseline ($S1$), we trained the model with the displacement reconstruction loss in (6) only. In $S2$, we added the standard $L1$ loss to $S1$, which corresponds to our loss in (7) without the landmark distance weights. To showcase the importance of weighting the contribution of each vertex, in $S3$ we added the landmark distance weights to the $L_{pr}$ loss. Results are shown in Table 2, where the remarkable improvement of our proposed loss against the standard $L1$ turns out evidently. This is explained by the fact that assigning a greater weight to movable face parts allows the network to focus on regions that are subject to strong facial motions, ultimately resulting in realistic samples.

### 4.3. 4D Facial Expression Evaluation

Since Motion3DGAN generates samples from noise to encourage diversity, the generated landmarks and meshes slightly change at each forward pass. Thus, computing the mean per-vertex error with respect to ground-truth shapes as done in [39] cannot represent a good and reliable measure

in this case. So, we evaluated the quality of the generated expression sequences implementing an expression classification solution, similar to [39]. We trained a classifier with one LSTM layer followed by a fully connected layer to recognize the 12 dynamic facial expressions of CoMA given a landmarks sequence as input. We trained this classifier on the same sequences used to train Motion3DGAN. The first expression sample of each identity form the test set, resulting in 144 testing samples.

Since neither the dataset nor the code in [39] are available for comparison, based on the information therein, we implemented a similar architecture relying on an LSTM to generate the per-frame expression and use it as baseline. The LSTM is trained to generate the motion of landmarks from an input code indicating the temporal evolution of the expression from neutral to the apex phase. Also this model was trained on the same data used for Motion3DGAN.

For testing, we generated 144 sequences with both Motion3DGAN and the LSTM generator. The sequences are consistent with those of the Motion3DGAN test set described above. These are then used to generate their corresponding meshes with our S2D-Dec. In Table 3, we report results in terms of classification accuracy and Frechet Inception Distance (FID) [26]. The metrics are computed either using the landmark sequences directly generated by Motion3DGAN and LSTM (Gen-LM row), or those extracted from the generated meshes (Det-LM row). Given that both Motion3DGAN and LSTM act as landmark generators, we also report the results obtained using sequences coming from a "perfect" generator, that is the ground truth landmark sequences of the test set (GT lands. column). This represents a sort of upper bound for the classification accuracy. We note that the features used to compute the FID metric are extracted by the last fully connected layer of our trained classifier, which outputs 512 features per sequence.

In Table 3, we observe that, in all the cases, Motion3DGAN surpasses the accuracy of LSTM to a large extent, providing a clear evidence that the generated sequences better capture the expression dynamics. This is also supported by the lower FID, which indicates that the Motion3DGAN samples better approximate the ground truth motions. The same conclusion is drawn from the closer recognition rate of Motion3DGAN to that obtained with the ground truth sequences. Furthermore, the accuracy increases by first generating the corresponding meshes and then re-extracting the landmarks from them. This suggests the S2D-Dec is capable of maintaining the particular motion, which is also supported by the similar recognition rate obtained with ground truth landmarks (73%) and those detected on their corresponding meshes generated with S2D-Dec (73.61%).
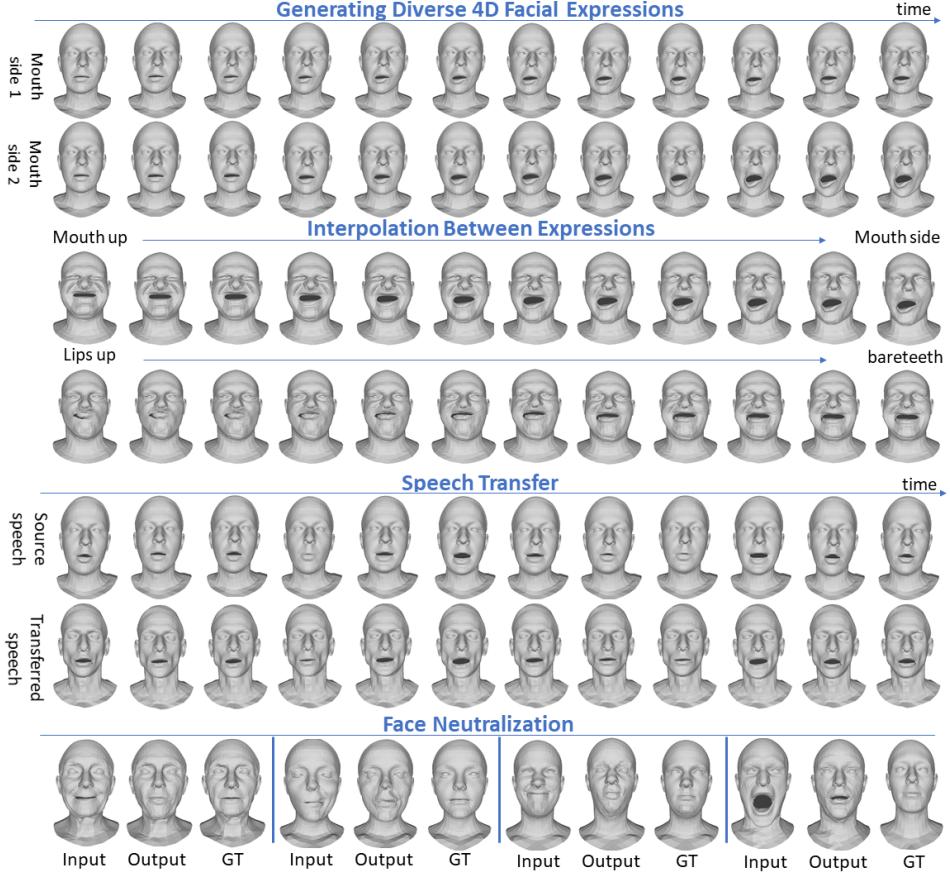
**Generating Diverse 4D Facial Expressions**    time

Mouth side 1

Mouth side 2

**Interpolation Between Expressions**

Mouth up      Mouth side

Lips up      bareteeth

**Speech Transfer**    time

Source speech

Transferred speech

**Face Neutralization**

Input   Output   GT    Input   Output   GT    Input   Output   GT    Input   Output   GT

Figure 5. **Applications** – From top to bottom: **Diversity**: the same identity performing the same class of facial expression (mouth side) with two different motions generated by Motion3DGAN. **Interpolation**: dynamic expressions resulting from the interpolation between two expressions peaks. **Transfer**: speech transfer from one identity to another. **Face neutralization**: for each of the four examples, we show the input expressive face, the neutralized face with S2D-Dec, and the ground truth neutral face of the given identity.

| Method | Classification Accuracy (%) ↑ | | | FID ↓ | |
|---|---|---|---|---|---|
| | GT lands | Mo3DGAN | LSTM | Mo3DGAN | LSTM |
| Gen-LM | **73.00** | **65.28** | 46.53 | **20.45** | 21.76 |
| Det-LM | **73.61** | **69.44** | 52.08 | **19.01** | 27.96 |

Table 3. Classification accuracy (%) and Frechet Inception Distance (FID) obtained with Ground Truth (GT) landmarks, Motion3DGAN and LSTM. Results are obtained using either the generated landmarks directly (Gen-LM), or by extracting landmarks from the meshes resulting from applying the S2D-Dec to the landmarks motion (Det-LM).

## 4.4. Applications

Our solution has some nice properties that open the way to various applications as shown in Figure 5.

**4D facial expression generation**: In the top two rows of Figure 5, we show the ability of Motion3DGAN to generate sequences for the same expression label that are highly variegated. In spite of that, S2D-Dec is able to generalize and reconstruct realistic meshes.

**Interpolation between facial expressions**: One interesting property of our Motion3DGAN is the possibility, enabled by the SRVF representation, of interpolating between generated motions. Given two points on the sphere $q_1$ and $q_2$, representing the motion sequences of two expressions, the geodesic path $\psi(\tau)$ between them is given by, $\psi(\tau) = \frac{1}{sin(\theta)}(sin(1-\tau)\theta)q_1 + sin(\theta\tau)q_2$, where, $\theta = d_{\mathcal{C}}(q_1, q_2) = \cos^{-1}(\langle q_1, q_2 \rangle)$. This path determines all the points $q_i$ existing between $q_1$ and $q_2$, each of them corresponding to a sequence of landmarks. Using our S2D-Dec, we can transform them to a 4D facial expression. Furthermore, while Motion3DGAN generates only neutral to apex sequences, we can exploit this interpolation to generate mixed 4D facial expressions that switch between the apex phase of different expressions by considering the last frame of each interpolated sequence. Figure 5 illustrates the interpolated faces between the apex frames of two expressions (examples are given for Mouth-up to Mouth-side, and

for Lips-up to Bare-teeth).

**Facial expression and speech transfer**: By using landmarks, our S2D-Dec can transfer facial expressions or speech between identities. This is done by extracting the sequence of landmarks from the source face, encoding their motion as an SRVF representation, transferring this motion to the neutral landmarks of the target face and using S2D-Dec to get the target identity following the motion of the first one. Some examples of speech transfer on the VO-CASET dataset [12] are shown in Figure 5 (see the supplementary material for their corresponding animations).

**Neutralization**: Given an input expressive face, S2D-Dec can generate the corresponding neutral face. This is obtained by introducing the displacements between the landmarks of an expressive face and those of a neutral template to S2D-Dec, so to generate the displacements needed to neutralize the expression. The last row of Figure 5 shows that our model can neutralize the expressions to a great extent, even though such motions do not occur at all in the training data.

## 5. Acknowledgments

## 6. Conclusions and Limitations

In this paper, we proposed a novel framework for dynamic 3D expression generation from an expression label, where two decoupled networks separately address modeling the motion dynamics and generating an expressive 3D face from a neutral one. We demonstrated the improvement with respect to previous solutions, and showed that using landmarks is effective in modeling the motion of expressions and the generation of 3D meshes. We also identified two main limitations: first, our S2D-Dec generates expression-specific deformations, and so cannot model identities. Moreover, while Motion3DGAN can generate diverse expressions and allows interpolating on the sphere to obtain complex facial expressions, the samples are of a fixed length (*i.e.*, 30 meshes, from neutral to apex). However, as shown in the applications, S2D-Dec can deal with motion of any length since it is independent from Motion3DGAN.

## 7. Appendix

### 7.1. Landmarks Configuration

In Figure 6 we show, for three different expressions, the configuration of landmarks used to guide the generation of the facial expression.
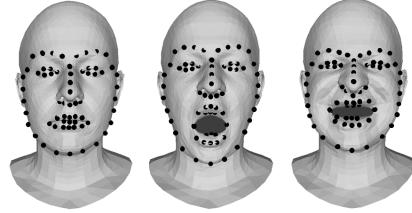


Figure 6. Landmarks configuration used to guide our model.

### 7.2. Logarithm and Exponential Maps

In order to map the SRVF data forth and back to a tangent space of $\mathcal{C}$, we use the logarithm $\log_p(.)$ and the exponential $\exp_p(.)$ maps defined in a given point $p$ by,

$$
\begin{aligned}
\log_p(q) &= \frac{d_{\mathcal{C}}(q,p)}{sin(d_{\mathcal{C}}(q,p))}(q - cos(d_{\mathcal{C}}(q,p))p), \\
\exp_p(s) &= cos(\|s\|)p + sin(\|s\|)\frac{s}{\|s\|},
\end{aligned}
\tag{8}
$$

where $d_{\mathcal{C}}(q,p) = \cos^{-1}(\langle q,p \rangle)$ is the distance between $q$ and $p$ in $\mathcal{C}$.

### 7.3. Architecture of S2D-Dec

The architecture adopted for S2D-Dec is based on the architecture proposed in [4]. S2D-Dec takes as input the displacements of 68 landmarks illustrated in Figure 6. The architecture includes a fully connected layer of size 2688, five spiral convolution layers of 64, 32, 32, 16 and 3 filters. Each spiral convolution layer is followed by an up-sampling by a factor of 4.

## 8. Ablation Study

In this section, we report a visual comparison between reconstructions obtained with the standard L1 loss and our proposed weighted L1. Figure 7 clearly shows the effect of our introduced weighting scheme that allows for improved expression modeling.

## References

[1] Victoria Fernandez Abrevaya, Stefanie Wuhrer, and Edmond Boyer. Multilinear autoencoder for 3D face model learning. In *IEEE Winter Conf. on Applications of Computer Vision (WACV)*, pages 1–9, 2018. 2
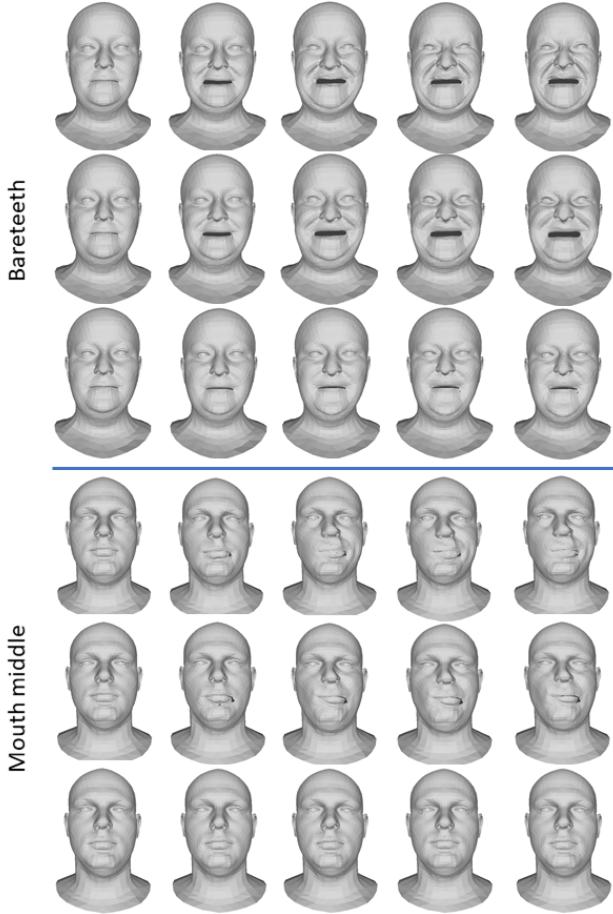
Figure 7. Ablation study: qualitative comparison between ground truth (first row) our model with (second row) and without (last row) weighted loss.

[2] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3D faces. In *Annual Conf. on Computer Graphics and Interactive Techniques (SIGGRAPH)*, pages 187–194, 1999. 2

[3] James Booth, Anastasios Roussos, Stefanos Zafeiriou, Allan Ponniah, and David Dunaway. A 3D morphable model learnt from 10,000 faces. In *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 5543–5552, 2016. 2

[4] Giorgos Bouritsas, Sergiy Bokhnyak, Stylianos Ploumpis, Stefanos Zafeiriou, and Michael Bronstein. Neural 3D morphable models: Spiral convolutional networks for 3D shape representation learning and generation. In *IEEE/CVF Int. Conf. on Computer Vision (ICCV)*, pages 7212–7221, 2019. 2, 5, 6, 9

[5] Michael M. Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst. Geometric deep learning: Going beyond euclidean data. *IEEE Signal Processing Mag.*, 34(4):18–42, 2017. 2

[6] Alan Brunton, Timo Bolkart, and Stefanie Wuhrer. Multilinear wavelets: A statistical shape space for human faces. In *European Conf. on Computer Vision*, pages 297–312. Springer, 2014. 2

[7] Chen Cao, Qiming Hou, and Kun Zhou. Displaced dynamic expression regression for real-time facial tracking and animation. *ACM Trans. on Graphics*, 33(4), July 2014. 2, 3

[8] Feng-Ju Chang, Anh Tuan Tran, Tal Hassner, Iacopo Masi, Ram Nevatia, and Gerard Medioni. Expnet: Landmark-free, deep, 3d facial expressions. In *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*, pages 122–129. IEEE, 2018. 3

[9] Feng-Ju Chang, Anh Tuan Tran, Tal Hassner, Iacopo Masi, Ram Nevatia, and Gerard Medioni. Faceposenet: Making a case for landmark-free face alignment. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 1599–1608, 2017. 3

[10] Lisha Chen, Hui Su, and Qiang Ji. Deep structured prediction for facial landmark detection. In *Advances in Neural Information Processing Systems (Neurips)*, volume 32, 2019. 3

[11] Darren Cosker, Eva Krumhuber, and Adrian Hilton. A facs valid 3D dynamic action unit database with applications to 3D dynamic morphable facial modeling. In *IEEE Int. Conf. on Computer Vision*, pages 2296–2303. IEEE, 2011. 5

[12] Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael Black. Capture, learning, and synthesis of 3D speaking styles. *Computer Vision and Pattern Recognition (CVPR)*, pages 10101–10111, 2019. 9

[13] Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael J. Black. Capture, learning, and synthesis of 3D speaking styles. In *IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 10093–10103, 2019. 2, 3

[14] Maxime Devanne, Hazem Wannous, Stefano Berretti, Pietro Pala, Mohamed Daoudi, and Alberto Del Bimbo. 3-D human action recognition by shape analysis of motion trajectories on Riemannian manifold. *IEEE Trans. on Cybernetics*, 45(7):1340–1352, 2014. 3

[15] Xuanyi Dong, Yi Yang, Shih-En Wei, Xinshuo Weng, Yaser Sheikh, and Shoou-I Yu. Supervision by registration and triangulation for landmark detection. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, pages 1–1, 2020. 3

[16] Hassen Drira, Boulbaba Ben Amor, Anuj Srivastava, Mohamed Daoudi, and Rim Slama. 3D face recognition under expressions, occlusions, and pose variations. *IEE Trans. on Pattern Analysis and Machine Intelligence*, 35(9):2270–2283, 2013. 3

[17] Lijie Fan, Wenbing Huang, Chuang Gan, Junzhou Huang, and Boqing Gong. Controllable image-to-video translation: A case study on facial expression generation. In *Conf. on Artificial Intelligence (AAAI) Symposium on Educational Advances in Artificial Intelligence*, pages 3510–3517. AAAI Press, 2019. 1

[18] Claudio Ferrari, Stefano Berretti, Pietro Pala, and Alberto Del Bimbo. Rendering realistic subject-dependent expression images by learning 3dmm deformation coefficients. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018. 3

[19] Claudio Ferrari, Stefano Berretti, Pietro Pala, and Alberto Del Bimbo. A sparse and locally coherent morphable face model for dense semantic correspondence across heterogeneous 3D faces. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 2021. 2, 6

[20] Claudio Ferrari, Giuseppe Lisanti, Stefano Berretti, and Alberto Del Bimbo. A dictionary learning-based 3D morphable shape model. *IEEE Trans. on Multimedia*, 19(12):2666–2679, 2017. 6

[21] Claudio Ferrari, Giuseppe Lisanti, Stefano Berretti, and Alberto Del Bimbo. Dictionary learning based 3D morphable model construction for face recognition with varying expression and pose. In *2015 International Conference on 3D Vision*, pages 509–517. IEEE, 2015. 2, 3

[22] Pablo Garrido, Levi Valgaerts, Ole Rehmsen, Thorsten Thormahlen, Patrick Perez, and Christian Theobalt. Automatic face reenactment. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4217–4224, 2014. 3

[23] Baris Gecer, Stylianos Ploumpis, Irene Kotsia, and Stefanos Zafeiriou. Ganfit: Generative adversarial network fitting for high fidelity 3d face reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1155–1164, 2019. 3

[24] Syed Zulqarnain Gilani, Ajmal Mian, and Peter Eastwood. Deep, dense and accurate 3d face correspondence for generating population specific deformable models. *Pattern Recognition*, 69:238–250, 2017. 3

[25] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *Advances in Neural Information Processing Systems (NIPS)*, pages 5767–5777, 2017. 4

[26] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 7

[27] Zhiwu Huang, Jiqing Wu, and Luc Van Gool. Manifold-valued image generation with wasserstein generative adversarial nets. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI*, pages 3886–3893. AAAI Press, 2019. 4

[28] Anis Kacem, Mohamed Daoudi, Boulbaba Ben Amor, and Juan Carlos Alvarez-Paiva. A novel space-time representation on the positive semidefinite cone for facial expression recognition. In *IEEE Int. Conf. on Computer Vision (ICCV)*, pages 3180–3189, 2017. 3

[29] Tero Karras, Timo Aila, Samuli Laine, Antti Herva, and Jaakko Lehtinen. Audio-driven facial animation by joint end-to-end learning of pose and emotion. *ACM Trans. on Graphics*, 36(4), July 2017. 2, 3

[30] Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4d scans. *ACM Trans. Graph.*, 36(6):194–1, 2017. 5

[31] Tianye Li, Timo Bolkart, Michael Julian, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4D scans. *ACM Trans. on Graphics, (Proc. SIGGRAPH Asia)*, 36(6), 2017. 3, 6

[32] Isaak Lim, Alexander Dielen, Marcel Campen, and Leif Kobbelt. A simple approach to intrinsic correspondence learning on unstructured 3D meshes. In *European Conf. on Computer Vision (ECCV) Workshops*, September 2018. 2

[33] Or Litany, Alex Bronstein, Michael Bronstein, and Ameesh Makadia. Deformable shape completion with graph convolutional autoencoders. In *IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, pages 1886–1895, 2018. 2

[34] Marcel Lüthi, Thomas Gerig, Christoph Jud, and Thomas Vetter. Gaussian process morphable models. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 40(8):1860–1873, 2017. 2

[35] Stylianos Moschoglou, Stylianos Ploumpis, Mihalis A. Nicolaou, Athanasios Papaioannou, and Stefanos Zafeiriou. 3dfacegan: Adversarial nets for 3D face representation, generation, and translation. *Int. Journal of Computer Vision*, 128:2534–2551, 2020. 2

[36] Thomas Neumann, Kiran Varanasi, Stephan Wenger, Markus Wacker, Marcus Magnor, and Christian

Theobalt. Sparse localized deformation components. *ACM Trans. on Graphics (TOG)*, 32(6):1–10, 2013. 2

[37] Naima Otberdout, Mohamed Daoudi, Anis Kacem, Lahoucine Ballihi, and Stefano Berretti. Dynamic facial expression generation on Hilbert hypersphere with conditional Wasserstein generative adversarial nets. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, pages 1–1, 2020. 1, 2, 3, 4

[38] Pascal Paysan, Reinhard Knothe, Brian Amberg, Sami Romdhani, and Thomas Vetter. A 3D face model for pose and illumination invariant face recognition. In *IEEE Int. Conf. on Advanced Video and Signal Based Surveillance*, pages 296–301, 2009. 2

[39] Rolandos Alexandros Potamias, Jiali Zheng, Stylianos Ploumpis, Giorgos Bouritsas, Evangelos Ververas, and Stefanos Zafeiriou. Learning to generate customized dynamic 3D facial expressions. In *European Conf. on Computer Vision (ECCV)*, pages 278–294, 2020. 3, 6, 7

[40] Anurag Ranjan, Timo Bolkart, Soubhik Sanyal, and Michael J. Black. Generating 3D faces using convolutional mesh autoencoders. In *European Conf. on Computer Vision (ECCV)*, pages 725–741, 2018. 2, 5, 6

[41] Anuj Srivastava, Eric Klassen, Shantanu H. Joshi, and Ian H. Jermyn. Shape analysis of elastic curves in euclidean spaces. *IEE Trans. on Pattern Analysis and Machine Intelligence*, 33(7):1415–1428, 2011. 3

[42] Nitika Verma, Edmond Boyer, and Jakob Verbeek. Feastnet: Feature-steered graph convolutions for 3D shape analysis. In *IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, pages 2598–2606, 2018. 2

[43] Jun Wan, Zhihui Lai, Jing Li, Jie Zhou, and Can Gao. Robust facial landmark detection by multiorder multiconstraint deep networks. *IEEE Trans. on Neural Networks and Learning Systems*, pages 1–14, 2021. 3

[44] Wei Wang, Xavier Alameda-Pineda, Dan Xu, Pascal Fua, Elisa Ricci, and Nicu Sebe. Every smile is unique: Landmark-guided diverse smile generation. In *IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, pages 7083–7092, 2018. 3

[45] Dan Zeng, Han Liu, Hui Lin, and Shiming Ge. Talking face generation with expression-tailored generative adversarial network. In *ACM Int. Conf. on Multimedia (MM'20)*, page 1716–1724, 2020. 3

[46] Xiangyu Zhu, Xiaoming Liu, Zhen Lei, and Stan Z Li. Face alignment in full pose range: A 3d total solution. *IEEE transactions on pattern analysis and machine intelligence*, 41(1):78–92, 2017. 3