

Locally Adaptive Neural 3D Morphable Models

Michail Tarasiou Rolandos Alexandros Potamias Eimear O’Sullivan
Stylianos Ploumpis Stefanos Zafeiriou
Imperial College London

{michail.tarasiou10,r.potamias19,e.o-sullivan,s.ploumpis,s.zafeiriou}@imperial.ac.uk

Abstract

We present the *Locally Adaptive Morphable Model (LAMM)*, a highly flexible Auto-Encoder (AE) framework for learning to generate and manipulate 3D meshes. We train our architecture following a simple self-supervised training scheme in which input displacements over a set of sparse control vertices are used to overwrite the encoded geometry in order to transform one training sample into another. During inference, our model produces a dense output that adheres locally to the specified sparse geometry while maintaining the overall appearance of the encoded object. This approach results in state-of-the-art performance in both disentangling manipulated geometry and 3D mesh reconstruction. To the best of our knowledge LAMM is the first end-to-end framework that enables direct local control of 3D vertex geometry in a single forward pass. A very efficient computational graph allows our network to train with only a fraction of the memory required by previous methods and run faster during inference, generating 12k vertex meshes at >60fps on a single CPU thread. We further leverage local geometry control as a primitive for higher level editing operations and present a set of derivative capabilities such as swapping and sampling object parts. Code and pretrained models can be found at <https://github.com/michaeltrs/LAMM>.

1. Introduction

The capacity to generate and manipulate digital 3D objects lies at the core of a multitude of applications related to the entertainment and media industries [9, 42], virtual and augmented reality [12, 24, 28, 35, 40] and healthcare [10, 23]. In particular, the ability to manually shape virtual humans can result in highly realistic avatars, while granting ultimate creative control to individual users.

However, achieving fine control in mesh manipulation necessitates the learning of a disentangled representation of 3D shapes which is still an open research problem [16–18]. Recently proposed methods based on Graph Convolutional

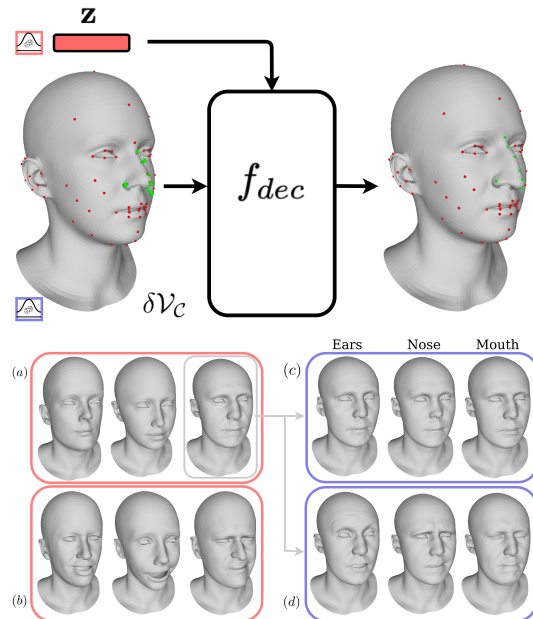


Figure 1. Overview of the *Locally Adaptive Morphable Model (LAMM)* use during inference. (top) Our trained decoder f_{dec} receives as inputs a latent code \mathbf{z} and displacements $\delta\mathcal{V}_c$ over a set of sparse control points (red vertices). Here displacements for control points in the nose region are shown with green arrows. The decoder generates the shape of the object encoded in \mathbf{z} , overwriting local geometry to respect $\delta\mathcal{V}_c$. (bottom) We can sample latent codes to generate new instances of human heads in the (a) identity, (b) expression space. Similarly, we can either randomly sample or provide vertex displacements manually at control points to manipulate (c) local identity features (ears, nose, mouth shown here) and (d) add expressions while retaining identity features unchanged.

Network (GCN)-based Auto-Encoders (AEs) [21, 26, 34] have demonstrated impressive performance in dimensionality reduction but typically learn a highly entangled latent space making them unsuitable for detailed shape manipulation. Additionally, despite having a low parameter count, these methods struggle to handle high-resolution meshes, limiting their applicability. Few works [17, 18] have dealt

with the disentanglement of local identity attributes, however these methods still rely on GCNs and opt for controlling manipulations through the state of the latent code which is partitioned and assigned to predefined object regions. Using the latent code to drive shape manipulation requires the use of explicit optimization objectives to learn a disentangled latent space. Moreover, partitioning its state is critically suboptimal for learning compressed representations of 3D objects. We propose a different paradigm which does not involve partitioning the latent code or relying on its state to drive changes in shape, resulting in state-of-the-art (SOTA) disentanglement and reconstruction capabilities in a unified architecture. Instead, we use a global latent code for 3D object unconditional generation and utilise additional inputs to jointly train our generative model to locally overwrite the latent encoded geometry. Our main contributions are the following:

- We present the *Locally Adaptive Morphable Model* (LAMM), a general framework for manipulating the geometry of registered meshes. To the best of our knowledge, this is the first method that allows direct shape control with a single forward pass. Applied on human 3D heads, LAMM exhibits SOTA disentanglement properties and allows for very fine geometric control over both facial identities and expressions.
- Our models, trained for manipulation, concurrently exhibit SOTA performance in mesh dimensionality reduction compared against methods trained exclusively on this task. As a result, a single model can be used to generate entirely new shapes and apply both localized and global modifications to their geometry.
- We show how our framework can leverage direct control as a primitive to achieve higher level editing operations such as region swapping and sampling.
- By deviating from GCN-based AE design, LAMM can scale to much larger meshes, needs only a fraction of GPU memory to train and can be significantly faster during inference compared to competing methods. For example, trained on 72k vertex meshes with batch size 32, our model requires 7.5Gb of GPU memory and runs at 0.045s on a single CPU thread. This model outperforms SpiralNet++ [21] which equivalently requires >40Gb of memory and runs $\times 13$ slower at 0.58s.
- We release training and evaluation codes to support future research. Additionally we will release pre-trained models and a Blender add-on [11] to facilitate intuitive mesh manipulation and fine-grained control.

2. Related Work

3D Morphable Models. Blanz and Vetter pioneered an innovative approach for synthesis and reconstruction of 3D human faces through their seminal *3D Morphable Model* (3DMM) [29]. First, they establish a fixed mesh topol-

ogy, ensuring that each face vertex is semantically significant and could be uniquely identified, e.g. the "tip of the nose". Subsequently, Principal Component Analysis (PCA) is employed to map 3D shapes and textures into a lower dimensional latent code. By adjusting these codes the PCA-based 3DMM can seamlessly generate new facial instances. Such 3DMMs are simple and very robust in terms of performance. They are well established as exceptionally reliable frameworks for 3D face/body analysis, finding applications in numerous human centric models [5, 16, 19, 24, 25, 29–32] which remain relevant to this date.

However, the linear formulation of PCA limits its expressivity and capacity to capture high frequency details. In an effort to capture nonlinear variations in the shape of human faces performing expressions [34] were the first to utilise Graph Convolutional Networks (GCN) for 3D dimensionality reduction. Their *Convolutional Mesh AutoEncoder* (COMA) makes use of fast and efficient spectral convolutions applied on a mesh surface [13] to compress and generate 3D meshes. An inductive bias more suitable for registered meshes is proposed in [7] who incorporate the spiral operator [26] in a GCN-based AE. This operator uses a spiral sequence to define an explicit local ordering for aggregating node features, thus giving the convolution kernels a sense of orientation across the mesh surface and greater expressive power. Following a similar principle, [21] further improve the performance of this architecture by enabling context aggregation across larger areas through dilated spiral operators. The deep GCN-based family of models discussed above have been consistently shown to outperform PCA at high levels of compression. However, when increasing the size of the latent space to values that are practical for applications, e.g. monocular 3D face reconstruction [6, 19], PCA is shown to match and quickly surpass GCN architectures. We argue that the translation equivariant property of GCNs is not optimal for modelling registered meshes as it does not account for the distinct semantics of each vertex. In contrast, by employing fully connected layers to encode input shapes into tokens and a global receptive field for downstream processing, our framework respects this attribute of the data, having the capacity to treat similar features differently depending on their location.

Face Manipulation. All statistical models of 3D meshes discussed above can be used to manipulate mesh geometry, e.g. the optimal fit for a given shape can be discerned through latent code optimization. The main issue with this approach is that latent spaces learnt for generative modelling are typically highly entangled [3, 16, 17]. This attribute discards the possibility for finely-tuned control of local geometry as varying any one dimension of the latent code can concurrently impact multiple properties of the reconstructed data. Numerous studies have been proposed to specifically address the challenge of en-

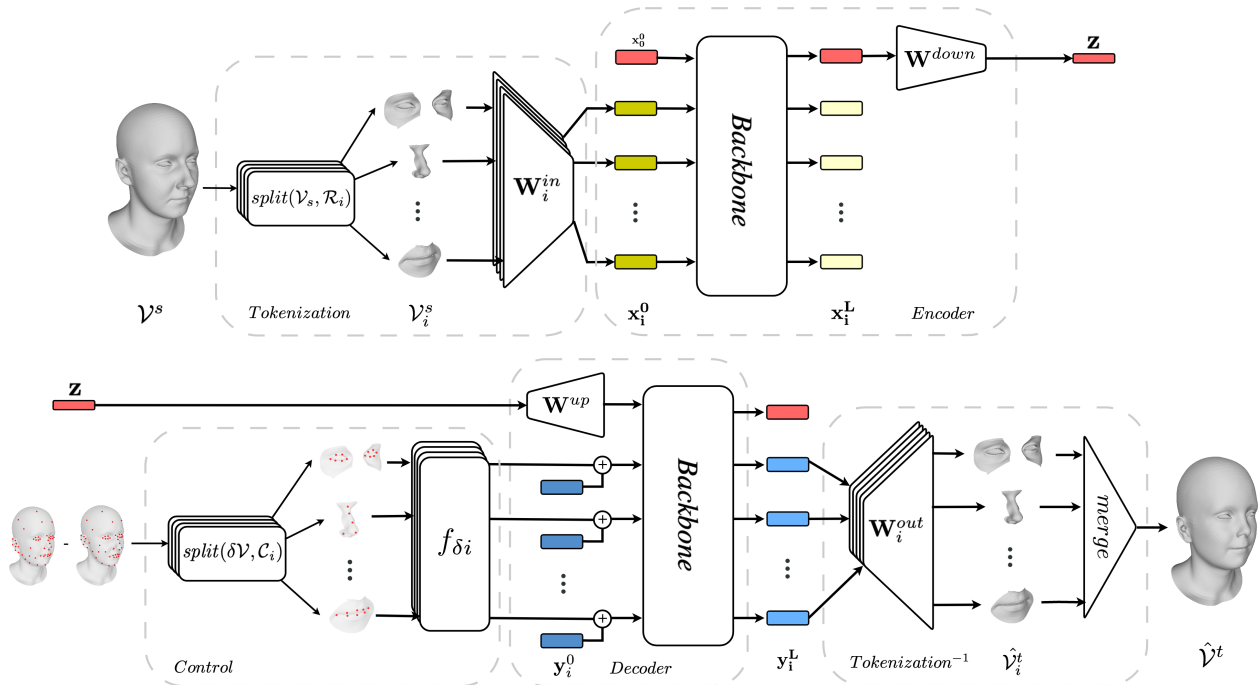


Figure 2. Architecture overview. A *source* mesh \mathcal{V}^s is encoded into latent code \mathbf{z} and decoded using additional displacements at control points $\delta\mathcal{V}_C$. (top) Tokenization and encoder modules. A 3D input mesh \mathcal{V}^s is split into regions $\mathcal{V}_i^s = \mathcal{V}^s[\mathcal{R}_i]$, each region is tokenized via region-specific linear weights \mathbf{W}_i^{in} and compressed by the encoder module into latent code \mathbf{z} . (bottom) Displacement control, decoder and inverse tokenization modules. User displacements at control points are split into corresponding regions \mathcal{C}_i , processed by region specific control networks f_{δ_i} and used by the decoder to overwrite the geometry encoded in \mathbf{z} . The inverse tokenization module translates decoder outputs into region geometries, via region-specific linear weights \mathbf{W}_i^{out} , which are merged together, using \mathcal{R}_i , into the *target* estimate $\hat{\mathcal{V}}^t$.

tangled representations of 3D face data for shape manipulation [16]. They predominantly utilize face regions to disentangle a latent representation and achieve local control through optimization, either by employing linear models [20] or GCN [2, 17, 18, 41]. Notably, works by [17, 18] adeptly deform 3D face geometry to manipulate an individual’s identity. Both works harness the backbone architecture from [21], along with various face regions, each governed by non-overlapping parts of the latent code. In [17] their AE learns a disentangled latent space by creating a composite batch through swapping arbitrary features across different samples. This enables the definition of a contrastive loss that leverages known differences and similarities in the latent representations. In a related approach, [18] also segment their latent codes to region specific parts and use a loss component forcing segments to adhere to the local eigenprojections of identity attributes. Similar to these works, we also employ predefined regions to segment our inputs. Importantly, we do not partition our latent codes into region-specific segments, which is beneficial for retaining strong performance in mesh reconstruction. Furthermore, we do not enforce spatial disentanglement through explicitly defined loss components (we train only with L_1 loss) but rather reach SOTA disentanglement performance solely

as an emerging property of our architecture.

3. Method

In this section we present the LAMM framework for the disentangled direct manipulation of 3D shapes. First, we describe the architecture, followed by a self-supervised training scheme for learning direct local control over mesh geometry. In the following, we assume a 3D mesh is represented as $\mathcal{M} = \{\mathcal{V}, \mathcal{F}\}$, where $\mathcal{V} \in \mathbb{R}^{N \times 3}$ is a set of N vertices represented as points in 3D space, and \mathcal{F} is a set of faces defining a shared topology.

3.1. Architecture

We leverage the AE framework for mesh representation learning, with targeted modifications in the decoder architecture to facilitate direct manipulation of mesh geometry. First, a *source* mesh \mathcal{V}^s is encoded into a latent space $\mathbf{z} = f_{enc}(\mathcal{V}^s) \in \mathbb{R}^D$ and subsequently decoded into an estimated *target* geometry $\hat{\mathcal{V}}^t = f_{dec}(\mathbf{z}, \delta\mathcal{V}_C)$ using control vertex displacements $\delta\mathcal{V}_C$ as additional input to the decoder. Both encoder and decoder modules utilise patch-based backbones, namely the Transformer [39] and MLP-Mixer [37], that operate on unordered sets of tokens and ex-

hibit a homogeneous feature space, i.e. feature dimensions stay constant throughout the network. An overview of the proposed architecture is presented schematically in Fig.2.

Tokenization. For deriving the inputs to the encoder we first split each *source* object into K non-overlapping dense regions, $\mathcal{V}_i^s = \mathcal{V}^s[\mathcal{R}_i] \in \mathbb{R}^{N_i \times 3}$, defined by a set of indices, $\mathcal{R}_i = \{j \in \mathbb{N}, 1 \leq j \leq N\}, i = 1, \dots, K$ (see Fig.3). We further flatten each region into a vector of size $\mathbf{v}_i^s \in \mathbb{R}^{3N_i}$ and tokenize it as a linear projection $\mathbf{W}_i^{in} \mathbf{v}_i^s$, where $\mathbf{W}_i^{in} \in \mathbb{R}^{D \times 3N_i}$. The proposed tokenization scheme differs significantly from the one used in [15, 37] that operate on fixed-size square patches. First, no parameter sharing takes place during this step. Instead there exists a one-to-one mapping between tokenization weights \mathbf{W}_i^{in} and input vertex coordinates that respects the semantic meaning of each vertex. This allows us to define regions with varying number of vertices that are designed to include distinct object parts based on the template mesh, e.g. a "whole eye" in human face meshes. Finally, our regions are typically much larger than commonly employed in computer vision [15, 36, 37], which results in a small number of tokens. For example, while the original implementation of ViT [15] used 256 tokens for image classification, we use as few as 11 tokens for 3D face modelling. Since the Transformer space and time complexities are both quadratic to the number of tokens, our models are very efficient in terms of memory consumption, can scale to very large meshes and can run significantly faster than GCNs especially on a CPU.

Encoder. Our fixed size tokenized representation for region i is $\mathbf{x}_i^0 = \mathbf{W}_i^{in} \mathbf{v}_i^s$. To these features we prepend a learnable identity token $\mathbf{x}_0^0 \in \mathbb{R}^D$ [14, 38] which is independent of the input. After processing by the L -layer homogeneous encoder, we retain only the state of the first index token $\mathbf{x}_0^L \in \mathbb{R}^D$ that forms the encoded representation of the *source* geometry. To obtain our latent code we linearly project this vector into a space of desired dimensionality $\mathbf{z} = \mathbf{W}^{down} \mathbf{x}_0^L \in \mathbb{R}^d$. The tokenization and encoder modules are illustrated in the top part of Fig.2.

Decoder. In the decoder module, we first linearly project the latent code back into a D dimensional vector $\mathbf{y}_0^0 = \mathbf{W}^{up} \mathbf{z}$. Analogously to the encoder, we build our decoder inputs by appending K additional learned tokens $\mathbf{y}_i^0 \in \mathbb{R}^D$ to the projected latent code \mathbf{y}_0^0 [8, 36]. We then select a set of non-overlapping, sparse vertex indices $\mathcal{C}_i \subseteq \mathcal{R}_i$ that are used as control points (Fig.3). We gather all desired displacements at control points per region $\delta \mathcal{V}_{\mathcal{C}_i} = (\mathcal{V}^t - \mathcal{V}^s)[\mathcal{C}_i] \in \mathbb{R}^{3|\mathcal{C}_i|}$ and process each with a fully connected network f_{δ_i} the output of which is added to respective learned tokens $\mathbf{y}_i^0, i \geq 1$. To naturally make our method work as an AE, we design f_{δ_i} without any bias terms. In this manner, for the special case of $\delta \mathcal{V}_{\mathcal{C}_i} = 0$ we get $f_{\delta_i}(0) = 0$ and LAMM reproduces the behaviour of an AE.

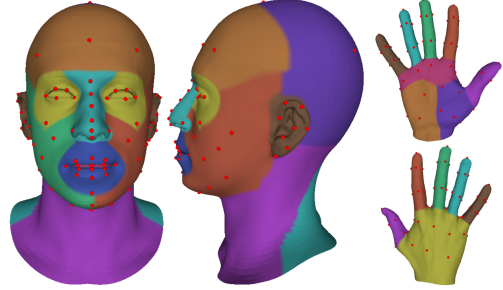


Figure 3. Templates for UHM12k and Handy data. Colors indicate dense regions \mathcal{R}_i . Control points \mathcal{C}_i are shown in red dots.

Tokenization⁻¹. Following the inverse process of region splitting and tokenization steps we now project each token independently to a vector $\hat{\mathbf{v}}_i^t = \mathbf{W}_i^{out} \mathbf{y}_i^L \in \mathbb{R}^{3N_i}$, reshape its dimensions to $N_i \times 3$ and use \mathcal{R}_i to merge regions into a unified 3D template mesh, \mathcal{V}^t . The decoder and inverse tokenization modules are presented schematically in the bottom part of Fig.2.

3.2. Training

Self-supervised objective. At each step we sample B *source* and *target* elements $\mathbf{V}^s, \mathbf{V}^t \in \mathbb{R}^{B \times N \times 3}$ from the training set and construct our *source* and *target* batches respectively as $\{\mathbf{V}^s, \mathbf{V}^s\}$ and $\{\mathbf{V}^s, \mathbf{V}^t\}$. In this manner, the first half of our batch is used for AE training and the second half for learning *source-to-target* transformations. For identity manipulation we use a random permutation of \mathbf{V}^s as \mathbf{V}^t . For expression manipulation, \mathbf{V}^s and \mathbf{V}^t include the same identities with neutral and non-neutral expressions.

Training with direct supervision from real *target* data has the benefit of providing realistic examples for the *target* geometry space. However, there are a few issues with this approach. In practice, interactive vertex control is an iterative process. As a result typical user input displacements are expected to be smaller and significantly sparser than the ones that globally transform one sample into another. To reduce the effect of displacement size discrepancy we train with a linear combination of *source* and *target* samples $\alpha \mathcal{V}^s + (1 - \alpha) \mathcal{V}^t, \alpha \sim \mathcal{U}(\alpha_{min}, 1)$ as the new *target* values. This effectively reduces the absolute values of displacements encountered during training in addition to augmenting the *target* space. The sparsity discrepancy is a trickier problem as it is not trivial to find realistic training pairs with only localized differences. We find that this issue is solved to a satisfactory degree implicitly as the proposed architecture naturally propagates deformations only to non-zero-displacement regions. We believe this attribute to be an intrinsic property of our architecture. In support of this claim, we show in the supplementary material that only through AE training, i.e. $\delta \mathcal{V}_{\mathcal{C}_i} = 0$, corrupting the values of learnable decoder tokens $\mathbf{y}_i^0, i > 0$ influences only their

corresponding region geometries.

Multilayer Loss. To adequately train our architecture, a distance-based loss between the *target* and output geometries is applied which suits our fundamental learning objective. In all experiments we utilize the $L1$ distance metric, calculated as $\|\mathcal{V}^t - \mathcal{V}^l\|_1$. We introduce an extension to this loss, applied at every level of the model, which increases the encoder’s capacity to encode input-related information and guides the decoder transformation from *mean* to *target* shape. Both employed backbones [37, 39] produce homogeneous feature spaces of size D across all encoder and decoder layers. At every layer l we decode region features $\mathbf{x}_i^l, \mathbf{y}_i^l, i \geq 1$ using \mathbf{W}_i^{out} and merge regions together to obtain a multilayer output geometry.

The objective of the encoder is to obtain a compressed representation of the input in the form of the latent code \mathbf{z} . At input, we want our tokens to retain as much information about the input shape to begin with, so we add a loss component $\|\mathbf{W}_i^{out} \mathbf{x}_i^0 - \mathbf{v}_i^s\|_1$ that encourages the input tokens to decode close to the *source* geometry. Furthermore, we provide some guidance to our encoder to progressively discard input related information from the region features through a loss component $\|\mathbf{W}_i^{out} \mathbf{x}_i^L - \bar{\mathbf{v}}_i\|_1$ that encourages the final encoder region tokens to match the mean shape for each region $\bar{\mathbf{v}}_i$ over the training data.

Similarly, we expand this pattern for all intermediate layers in a linear fashion.

$$L_{enc} = \sum_{l=0}^L \frac{\lambda^l}{K} \sum_{i=1}^K \|\mathbf{W}_i^{out} \mathbf{x}_i^l - \frac{1}{L}((L-l)\mathbf{v}_i^s + l\bar{\mathbf{v}}_i)\|_1 \quad (1)$$

In the decoder, we first center initialised tokens \mathbf{y}_i^0 through the term $\|\mathbf{W}_i^{out} \mathbf{y}_i^0 - \bar{\mathbf{v}}_i\|_1$ which encourages them to decode into the mean shape. The output of the final layer will be our model’s estimate of the object geometry, thus, we supervise through the term $\|\mathbf{W}_i^{out} \mathbf{y}_i^L - \mathbf{v}_i^t\|_1$. Similarly, we apply a multilayer loss to intermediate layers to encourage a progressive transformation of the region tokens from *mean* to *target* geometry.

$$L_{dec} = \sum_{l=0}^L \frac{\lambda^l}{K} \sum_{i=1}^K \|\mathbf{W}_i^{out} \mathbf{y}_i^l - (\frac{1}{L}(l\mathbf{v}_i^t + (L-l)\bar{\mathbf{v}}_i)\|_1 \quad (2)$$

In both eqs.1, 2 λ_l is a weight for the loss component at layer l . In practice we opt for a simple solution of keeping all $\lambda_l = 1$ which leads to stable training and good generalization performance. We note that in both equations the inverse tokenization weights \mathbf{W}_i^{out} are shared in every loss term, forcing a common geometric representation throughout the network which progressively transforms from *source* to *mean* to *target*. Additionally, our data are

typically mean centered prior to any processing, thus, $\bar{\mathbf{v}}_i$ reduces to the zero vector. In this case, our multilayer loss can be viewed as a form of regularization forcing the values of \mathbf{W}_i^{out} to be either small in values or orthogonal to the learnt tokens \mathbf{y}_i^0 .

4. Experiments

4.1. Implementation details

Datasets. For effectively learning disentangled human head identity manipulations we need a large scale dataset of 3D heads in neutral expression. For this reason, we register 6k heads from the original data used in the UHM model [31] into a new quads-based template which consists of 12k vertices, is more friendly towards rigging and is animation ready. We will refer to this data as UHM12k. To create our synthetic expression data we applied a set of 28 blend-shapes [1] to the UHM12k data. To test our method on high resolution meshes, we utilise the linear models provided by the UHM [31] and sample 10k identities consisting of 72k vertices. To evaluate our method on hand data, we utilized Handy [33], a large-scale hand model, trained from more than 1k distinct identities. In particular, to generate our hand dataset, we randomly sampled 10k hands from the Handy latent space. All datasets were split into training and evaluation sets at a 9-1 ratio.

Architecture. We optimize our architecture through an ablation study on autoencoding using the UHM12k data, presented in the supplementary material. Overall, our architecture consists of five encoder and three decoder layers with 512 feature dimension. For manipulation experiments we split the UHM12k and the Handy templates into 11 and 9 regions and utilise 73 (101 for expression modelling) and 72 control points respectively, all of which are illustrated in Figs.3, 5, 6. For identity control with the UHM data we use the same face regions as [17, 18] and 91 control points.

Training. Our proposed architecture can train for significantly longer than GCN baselines without observing a drop in evaluation set performance. We train all networks (including baselines) in mesh reconstruction and manipulation for 1,500 epochs on a single Nvidia Titan X GPU, with batch size 32 (adjusted for GCNs according to memory). We use the AdamW optimizer [22] starting with a 10 epoch learning rate warm up to a 10^{-4} which is decayed via a 1-cycle cosine decay [27] to 10^{-6} . For local control experiments, we find it is crucial to initialize our models from checkpoints pre-trained exclusively for autoencoding. Further details are provided in the supplementary material.

4.2. 3D shape reconstruction

In Table 1 we compare the reconstruction performance of our proposed framework for autoencoding against PCA [4] and state of the art GCNs [7, 21, 34]. Overall, our method

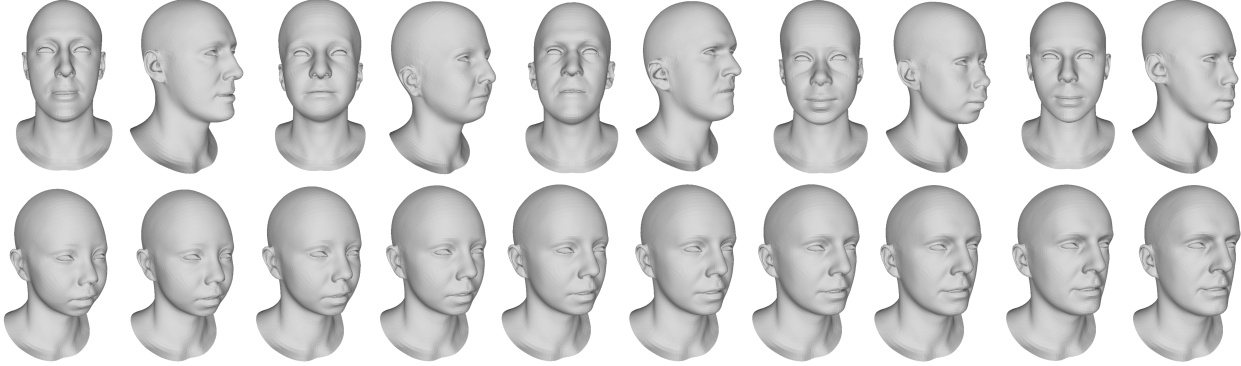


Figure 4. Qualitative assessment of trained AE in UHM12k. (top) New shape generation through sampling of the latent space \mathbf{z} , (bottom) Interpolation between latent codes for two samples from the evaluation data shows a smooth transition between identities.

Table 1. Quantitative evaluation of 3D shape reconstruction for models trained exclusively in autoencoding with latent size 256. Presented values are mean Euclidean distances ($\times 10^{-2}$ mm).

	UHM12k	+expr.	UHM	Handy
PCA	10.42	11.49	11.73	25.90
COMA [34]	13.11	14.20	15.40	27.20
Neural3DMM [7]	11.82	12.97	13.88	26.75
SpilarNet++ [21]	11.53	12.58	13.55	26.72
LAMM-Transformer	9.08	9.09	13.70	26.31
LAMM-MLPMixer	7.97	9.51	11.48	24.60

is shown to significantly outperform others applied in head and hand data with improvements being more pronounced on UHM12k.

Comparing the two employed backbones, the MLP Mixer is a consistently strong performer, surpassing baselines in all cases tested. Transformer-based backbones are found to be more suitable for the expression enhanced data, but weaker than the MLP Mixer for the UHM and Handy data. In Fig.4 (top) we present generated samples drawn from our head model. Only the decoder part of our architecture is used here. Similar to previous works [4, 7, 21, 34] sampling new instances of global shape can be achieved by fitting a Gaussian probability distribution to \mathbf{z} values collected over the training data.

In the bottom part of Fig.4 we choose two highly diverse identities from the evaluation set, encode both and present geometries generated by linear interpolation of their latent codes, noting a smooth transition between generated shapes.

4.3. Local control

Global manipulation. We present a quantitative evaluation of models trained in shape manipulation in Table 2. First, we note that AE performance is very similar to the values presented in Table 1 where models are trained exclusively in mesh reconstruction. In fact, AE performance is

found to improve in many cases suggesting the complementarity of the two tasks. We also find that expression training can increase AE performance significantly. If we train for expression manipulation but select our model based on its AE performance our Transformer based model achieves a distance error of 7.79×10^{-2} mm, which is 14% less compared to the best value we obtained from AE training in Table 1 (9.08×10^{-2} mm). Figs.5, 6 illustrate per-vertex mean prediction errors plotted over respective templates for identity and expression manipulations. In Fig.5 these are shown to be clearly smaller close to control points verifying the metrics reported in Table 2. For human heads, errors are mostly concentrated in the frontal region of the face. For hands, errors are significantly larger and more evenly spread. However, there are clear concentrations near the wrist which can be attributed to the absence of control points in these regions. We additionally show examples of *source* to *target* generations from the evaluation data, observing how predicted geometries match the appearance of *target* for both datasets. For expression manipulation, we can generally observe from Fig.6 and Table 2 that errors are smaller compared to identity manipulation. As expected, we observe that errors are also concentrated in the face region where most of the deformation takes place and particularly small in head regions that do not participate in facial expressions such as the ears and skull. We also note that, in contrast to identity manipulation, errors are significantly lower here at non-control points which can be attributed to the fact that $\notin \mathcal{C}_i$ includes many vertices with very small deformations, e.g. skull.

Region swap. For swaping regions among meshes it is critical to find a set of control point displacements that best transform the *source* region \mathcal{V}_i^s to a desired *target* shape \mathcal{V}_i^t . However, calculating displacements directly results in non realistic geometries as the two regions may be offset by a significant amount. To resolve this issue, we first rigidly align regions by their common outer edge indices. Due to

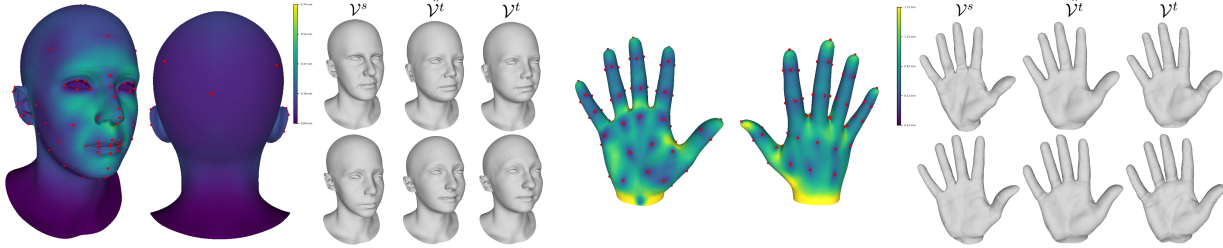


Figure 5. Identity manipulation performance under direct vertex control for UHM12k (left) and Handy (right) data. We observe that errors are typically reduced close to control points. For both datasets our method learns to produce highly realistic output geometries.

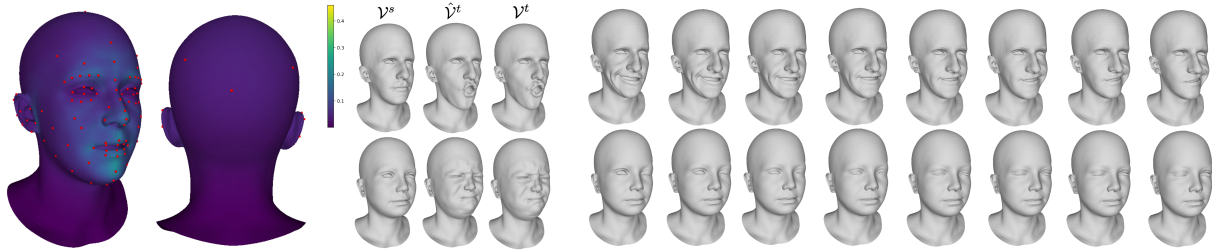


Figure 6. Expression manipulation performance under direct vertex control for UHM12k data. (left) We observe that per-vertex errors are very small and that \hat{V}^t closely resembles the appearance of V^t . (right) Linear interpolation in expression intensities for the two samples. We observe a smooth transition between facial expressions and the ability to perform highly localized edits by manipulating eyes independently.

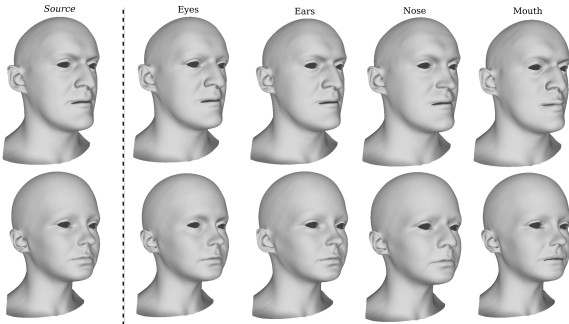


Figure 7. Head region swapping for UHM data.

Table 2. Quantitative evaluation of models trained in 3D shape manipulation. We present mean Euclidean distance error ($\times 10^{-2}$ mm) for the AE case ($V^t = V^s$), control ($\in C$) and remaining ($\notin C$) vertices. We note AE performance can surpass that of models trained exclusively for mesh reconstruction. † neutral expression only.

Backbone	UHM12k			UHM12k+expr.		
	AE	$\in C$	$\notin C$	AE†	$\in C$	$\notin C$
Transformer	8.79	18.47	26.03	8.61	15.86	9.67
MLPMixer	8.36	19.87	28.08	8.98	15.92	9.73
Backbone	UHM			Handy		
	AE	$\in C$	$\notin C$	AE	$\in C$	$\notin C$
Transformer	15.86	26.23	38.80	26.78	28.32	55.69
MLPMixer	11.63	17.97	32.03	24.59	25.57	51.66

the simplicity of applied transformations, rigid alignment

is found not to disturb the appearance of the *target* shape, keeping it realistic. The outcome of this operation is reduced offsets at the boundaries, however, there are still discontinuities due to the restricted nature of rigid alignment. It is critical to note, however, that the region’s control vertices, sparsely located in its interior, represent a valid sparse geometry since there exists a valid dense shape that fits these vertices while respecting the region’s boundary conditions. After finding the optimal rigid transformation, it can be applied to the control points of the *target* region to obtain the updated *target* locations of the control vertices, which in turn are used to calculate δV_{Ci} . We employ this methodology for region swapping in human heads. Fig.7 demonstrates swapping the eyes, ears, nose and mouth regions between two diverse meshes from the UHM data.

Region sampling and disentanglement. We show it is possible to generate new dense region shapes by sampling valid displacements for the control vertices. Using the same approach described for region swapping, we randomly select data pairs from the training data and collect a set of valid control vertex displacements for every region. To these, we fit a Gaussian distribution from which we can sample from to produce new valid control point displacements. We present samples of generated face regions in the top part of Fig.8 using the UHM data. In the bottom part, we utilise UHM12k to illustrate face part interpolation for the ears and nose regions. We evaluate the capacity of our model to have only localized effects in manipulated geometry in Fig.9 compared to SOTA method SD-VAE [17]. Con-

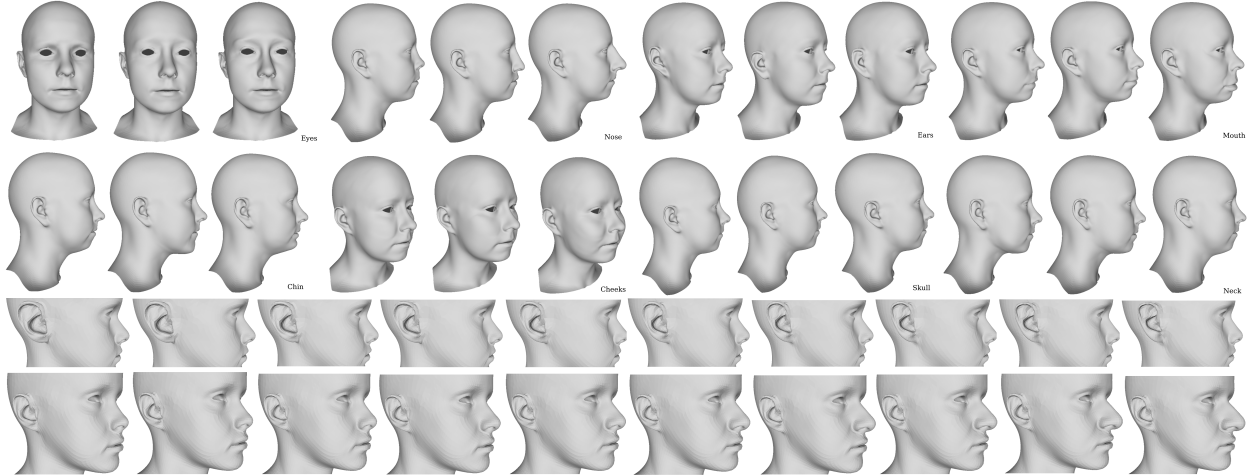


Figure 8. Random generation and interpolation of head regions.(top) Random head region generation on UHM (bottom) Region interpolation for ears and nose using UHM12k. We observe a smooth transition between generated regions while remaining areas are unaffected.

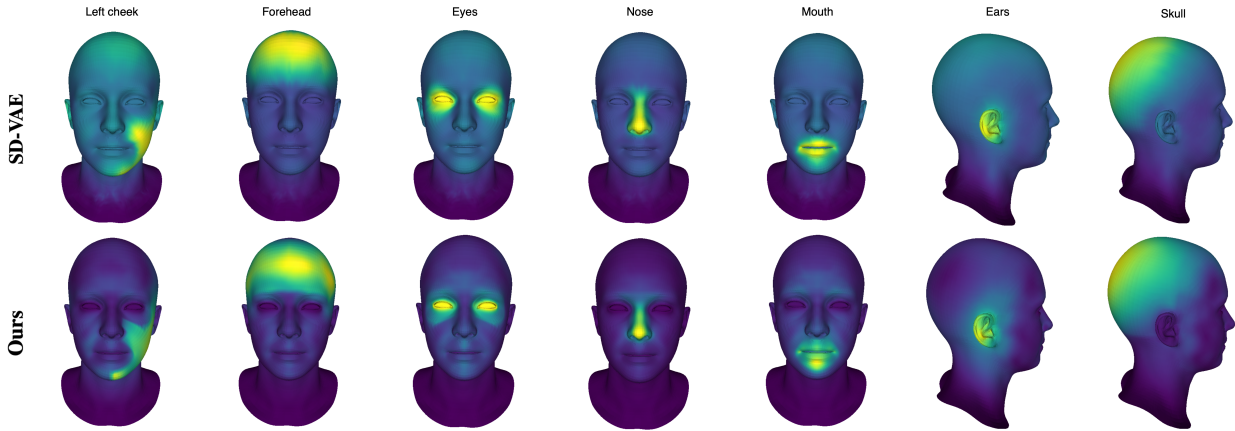


Figure 9. Evaluation of disentanglement performance compared to state-of-the-art SD-VAE [17]. For all data in the evaluation set we randomly sample 20 shapes for each head region. Plotted values are per-vertex mean Euclidean distances between input and output shapes.

sistently, we find that our method leads to more localized effects for all face regions and improved disentanglement over generated geometries. Importantly, we note that in contrast to [17], our framework is trained without any loss component aimed towards disentanglement which indicates this to be an intrinsic property of our method. We provide additional evidence to support this claim in the supplementary material. Furthermore, we find that SD-VAE has low reconstruction performance at $18.90 \times 10^{-2}mm$, more than double the reconstruction error of our method from Table 2 and more than SpiralNet++ [21] which is the backbone employed by SD-VAE. As we discuss in the supplementary material, this is mainly attributed to the choice of partitioning the latent code made by SD-VAE to control region geometries.

5. Conclusion

We have introduced LAMM, a comprehensive framework for learning 3D mesh representations and achieving spatially disentangled shape manipulation. To our knowledge, LAMM is the first end-to-end trainable method that is capable of directly manipulating mesh shape, by using desired displacements as inputs, in a single forward pass. Our extensive experiments with 3D human head and hand mesh data demonstrate that LAMM simultaneously attains state-of-the-art performance in both disentanglement and reconstruction within a unified architecture. Notably, it scales more effectively to high-resolution data, requires less memory for training, and offers faster inference speeds compared to existing methods. When implemented within the identity space of human heads, our method exhibits significant po-

tential for applications in face sculpting. Additionally, when trained in the expression space of human heads, LAMM demonstrates considerable promise for manually enhancing digital characters with expressive details and for its application in the field of animation. Fast CPU inference makes LAMM an energy efficient and economically viable option, democratizing access to advanced 3D modeling.

References

- [1] 3d scan store. <https://www.3dscanstore.com/>. Accessed: 2023-11-13. **5**
- [2] Mohammadamin Aliari, Tiberiu Beauchamp, AndrePopa, and Eric Paquette. Face editing using part-based optimization of the latent space. *Computer Graphics Forum*, 2023. **3**
- [3] Tristan Aumentado-Armstrong, Stavros Tsogkas, Allan Jepson, and Sven Dickinson. Geometric disentanglement for generative latent shape models. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 8180–8189, 2019. **2**
- [4] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Techniques*, page 187–194, USA, 1999. ACM Press/Addison-Wesley Publishing Co. **5, 6**
- [5] James Booth, Anastasios Roussos, Stefanos Zafeiriou, Allan Ponniah, and David Dunaway. A 3d morphable model learnt from 10,000 faces. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5543–5552, 2016. **2**
- [6] James Booth, Epameinondas Antonakos, Stylianos Ploumpis, George Trigeorgis, Yannis Panagakis, and Stefanos Zafeiriou. “3d face morphable models” in-the-wild”. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 48–57, 2017. **2**
- [7] Giorgos Bouritsas, Sergiy Bokhnyak, Stylianos Ploumpis, Michael Bronstein, and Stefanos Zafeiriou. Neural 3d morphable models: Spiral convolutional networks for 3d shape representation learning and generation. In *The IEEE International Conference on Computer Vision (ICCV)*, 2019. **2, 5, 6**
- [8] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Computer Vision – ECCV 2020*, pages 213–229, Cham, 2020. Springer International Publishing. **4**
- [9] Lucy Chai, Richard Tucker, Zhengqi Li, Phillip Isola, and Noah Snavely. Persistent nature: A generative model of unbounded 3d worlds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20863–20874, 2023. **1**
- [10] Richard J Chen, Ming Y Lu, Tiffany Y Chen, Drew F K Williamson, and Faisal Mahmood. Synthetic data in machine learning for medicine and healthcare. *Nat. Biomed. Eng.*, 5(6):493–497, 2021. **1**
- [11] Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. **2**
- [12] Enric Corona, Mihai Zanfir, Thiemo Alldieck, Eduard Gabriel Bazavan, Andrei Zanfir, and Cristian Sminchisescu. Structured 3d features for reconstructing relightable and animatable avatars. In *CVPR*, 2023. **1**
- [13] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2016. **2**
- [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, 2019. Association for Computational Linguistics. **4**
- [15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. **4**
- [16] Bernhard Egger, William A. P. Smith, Ayush Tewari, Stefanie Wuhler, Michael Zollhoefer, 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. *arXiv e-prints*, 2019. **1, 2, 3**
- [17] Simone Foti, Bongjin Koo, Danail Stoyanov, and Matthew J. Clarkson. 3d shape variational autoencoder latent disentanglement via mini-batch feature swapping for bodies and faces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18730–18739, 2022. **1, 2, 3, 5, 7, 8**
- [18] Simone Foti, Bongjin Koo, Danail Stoyanov, and Matthew J. Clarkson. 3d generative model latent disentanglement via local eigenprojection. *Computer Graphics Forum*, 42(6):e14793, 2023. **1, 3, 5**
- [19] 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. **2**
- [20] Donya Ghafourzadeh, Sahel Fallahdoust, Cyrus Rahgoshay, Andre Beauchamp, Adeline Aubame, Tiberiu Popa, and Eric Paquette. Local control editing paradigms for part-based 3D face morphable models. *Computer Animation and Virtual Worlds*, 32(6):e2028. **3**
- [21] Shunwang Gong, Lei Chen, Michael Bronstein, and Stefanos Zafeiriou. Spiralnet++: A fast and highly efficient mesh convolution operator. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019. **1, 2, 3, 5, 6, 8**

- [22] Loshchilov Ilya, Hutter Frank, et al. Decoupled weight decay regularization. *Proceedings of ICLR*, 2019. 5
- [23] Firas Khader, Gustav Müller-Franzes, Soroosh Tayebi Arasteh, Tianyu Han, Christoph Haarbuerger, Maximilian Schulze-Hagen, Philipp Schad, Sandy Engelhardt, Bettina Baeßler, Sebastian Foersch, Johannes Stegmaier, Christiane Kuhl, Sven Nebelung, Jakob Nikolas Kather, and Daniel Truhn. Denoising diffusion probabilistic models for 3D medical image generation. *Scientific Reports*, 13:7303, 2023. 1
- [24] Alexandros Lattas, Stylianos Moschoglou, Stylianos Ploumpis, Baris Gecer, Jiankang Deng, and Stefanos Zafeiriou. Fitme: Deep photorealistic 3d morphable model avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8629–8640, 2023. 1, 2
- [25] Tianye Li, Timo Bolkart, Michael J. Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4D scans. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 36(6):194:1–194:17, 2017. 2
- [26] Isaak Lim, Alexander Dielen, Marcel Campen, and Leif Kobbelt. A simple approach to intrinsic correspondence learning on unstructured 3d meshes. In *Computer Vision – ECCV 2018 Workshops: Munich, Germany, September 8-14, 2018, Proceedings, Part III*, page 349–362, Berlin, Heidelberg, 2019. Springer-Verlag. 1, 2
- [27] Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations*, 2017. 5
- [28] Foivos Paraperas Papanтониου, Alexandros Lattas, Stylianos Moschoglou, and Stefanos Zafeiriou. Relightify: Relightable 3d faces from a single image via diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 1
- [29] Pascal Paysan, Reinhard Knothe, Brian Amberg, Sami Romdhani, and Thomas Vetter. A 3d face model for pose and illumination invariant face recognition. In *2009 Sixth IEEE International Conference on Advanced Video and Signal Based Surveillance*, pages 296–301, 2009. 2
- [30] Stylianos Ploumpis, Haoyang Wang, Nick Pears, William AP Smith, and Stefanos Zafeiriou. Combining 3d morphable models: A large scale face-and-head model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10934–10943, 2019.
- [31] Stylianos Ploumpis, Evangelos Ververas, Eimear O’Sullivan, Stylianos Moschoglou, Haoyang Wang, Nick Pears, William AP Smith, Baris Gecer, and Stefanos Zafeiriou. Towards a complete 3d morphable model of the human head. *IEEE transactions on pattern analysis and machine intelligence*, 43(11):4142–4160, 2020. 5
- [32] Stylianos Ploumpis, Stylianos Moschoglou, Vasileios Triantafyllou, and Stefanos Zafeiriou. 3d human tongue reconstruction from single “in-the-wild” images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2771–2780, 2022. 2
- [33] Rolandos Alexandros Potamias, Stylianos Ploumpis, Stylianos Moschoglou, Vasileios Triantafyllou, and Stefanos Zafeiriou. Handy: Towards a high fidelity 3d hand shape and appearance model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4670–4680, 2023. 5
- [34] Anurag Ranjan, Timo Bolkart, Soubhik Sanyal, and Michael J. Black. Generating 3D faces using convolutional mesh autoencoders. In *European Conference on Computer Vision (ECCV)*, pages 725–741, 2018. 1, 2, 5, 6
- [35] Michael Stengel, Koki Nagano, Chao Liu, Matthew Chan, Alex Trevithick, Shalini De Mello, Jonghyun Kim, and David Luebke. Ai-mediated 3d video conferencing. In *ACM SIGGRAPH Emerging Technologies*, 2023. 1
- [36] Michail Tarasiou, Erik Chavez, and Stefanos Zafeiriou. Vits for sats: Vision transformers for satellite image time series. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10418–10428, 2023. 4
- [37] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, and Alexey Dosovitskiy. Mlp-mixer: An all-mlp architecture for vision. In *Advances in Neural Information Processing Systems*, pages 24261–24272. Curran Associates, Inc., 2021. 3, 4, 5
- [38] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve Jegou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357, 2021. 4
- [39] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc., 2017. 3, 5
- [40] Tengfei Wang, Bo Zhang, Ting Zhang, Shuyang Gu, Jianmin Bao, Tadas Baltrusaitis, Jingjing Shen, Dong Chen, Fang Wen, Qifeng Chen, and Baining Guo. Rodin: A generative model for sculpting 3d digital avatars using diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4563–4573, 2023. 1
- [41] Peizhi Yan, James Gregson, Qiang Tang, Rabab Ward, Zhan Xu, and Shan Du. Neo-3df: Novel editing-oriented 3d face creation and reconstruction. In *Proceedings of the Asian Conference on Computer Vision*, pages 486–502, 2022. 3
- [42] Haotian Zhang, Ye Yuan, Viktor Makoviyshuk, Yunrong Guo, Sanja Fidler, Xue Bin Peng, and Kayvon Fatahalian. Learning physically simulated tennis skills from broadcast videos. *ACM Trans. Graph.* 1