





Realistic One-shot Mesh-based Head Avatars

Taras Khakhulin^{1,2} , Vanessa Sklyarova^{1,2} ,
Victor Lempitsky³ , and Egor Zakharov^{1,2} 

¹ Samsung AI Center – Moscow

² Skolkovo Institute of Science and Technology

³ Yandex Armenia

<https://samsunglabs.github.io/rome/>

Abstract. We present a system for the creation of realistic one-shot mesh-based (ROME) human head avatars. From a single photograph, our system estimates the head mesh (with person-specific details in both the facial and non-facial head parts) as well as the neural texture encoding, local photometric and geometric details. The resulting avatars are rigged and can be rendered using a deep rendering network, which is trained alongside the mesh and texture estimators on a dataset of in-the-wild videos. In the experiments, we observe that our system performs competitively both in terms of head geometry recovery and the quality of renders, especially for cross-person reenactment.

1 Introduction

Personalized human avatars are becoming a key technology across several application domains, such as telepresence, virtual worlds, and online commerce. In many practical cases, it is sufficient to personalize only a part of the avatar’s body, while the remaining areas can then be picked from a pre-defined library of assets or omitted from the interface. Towards this end, many applications require personalization at the head level, i.e., the creation of person-specific head models, thus making it an important and viable intermediate step between personalizing only the face and the entire body. Alone, face personalization is often insufficient, while the full-body reconstruction remains a complicated task and leads to the reduced quality of the models or requires cumbersome data collection.

Acquiring human avatars from a single photograph (in a “one-shot” setting) offers the highest convenience for the end-user. However, their creation process is particularly challenging and requires strong priors on human geometry and appearance. To this end, parametric models are long known to offer good personalization solutions [3] and were recently adapted to one-shot performance [9,13,41]. Such models can be learned from a relatively small dataset of 3D scans and represent geometry and appearance via textured meshes, making them compatible with many computer graphics applications and pipelines. However, they cannot be trivially expanded to the whole head region due to the large geometric variability of the non-facial parts such as hair and neck.



Fig. 1: Our system creates realistic mesh-based avatars from a single **source** photo. These avatars are rigged, i.e., they can be driven by the animation parameters from a different **driving** frame. At the same time, our obtained **meshes** and **renderers** achieve a high degree of personalization in both appearance and geometry and are trained in an end-to-end fashion on a dataset of in-the-wild videos without any additional 3D supervision.

Our proposed system addresses this issue and allows parametric face models to represent the non-facial parts. In order to handle the increased geometric and photometric variability, we train our method on a large dataset of in-the-wild videos [6] and use neural networks to parameterize both the geometry and the appearance. For the appearance modeling, we follow the deferred neural rendering [46] paradigm and employ a combination of neural textures and rendering networks. In addition, a neural rendering framework [36] is used to enable end-to-end training and achieve high visual realism of the resulting head models. After training, the geometric and appearance networks can be conditioned on the information extracted from a single photograph, enabling one-shot realistic avatar generation.

To the best of our knowledge, our system is the first that is capable of creating realistic personalized human head models in a rigged mesh format from a single photograph. This distinguishes our model from a growing class of approaches that a) recover neural head avatars without explicit geometry [49,42,53,52], b) can personalize the face region but not the whole head [45,20,3,9], and c) from commercial systems that create non-photorealistic mesh avatars from a single image [1,33]. Alongside our main model, we also discuss its simplified version based on a linear blendshapes basis and show how to train it using the same video dataset. Below, we refer to the avatars generated by our system as ROME avatars (Realistic One-shot Mesh-based avatars).

2 Related work

Parametric models of human faces. Over the recent decades, 3D face reconstruction methods have been actively employed to tackle the problems of face tracking and alignment [15,13], face recognition [2,47], and generative modelling [45,20,29,35,25,26]. In all these scenarios, statistical mesh-based models, aka parametric models [3], remain one of the widely used tools [8,34]. State-of-the-art parametric models for human heads consist of rigged meshes [23] which

support a diverse range of animations via disentangled shape and expression blendshapes and rigid motions for the jaw, neck, and eyeballs. However, they only provide reconstructions for the face, ears, neck, and forehead regions, limiting the range of applications. Including full head reconstruction (i.e., hair and shoulders) into these parametric models is possible, but existing approaches require significantly more training data to be gathered in the form of 3D scans. Instead, in our work, we propose to leverage existing large-scale datasets [6] of in-the-wild videos via the learning-by-synthesis paradigm without any additional 3D annotations.

Neural 3D human head models. While parametric models provide sufficient reconstruction quality for many downstream applications, they are not able to depict the fine appearance details that are needed for photorealistic modelling. In recent years, the problem of representing complex geometry and appearance of humans started being addressed using high-capacity deep neural networks. Some of these works use strong human-specific priors [35,9,39,27], while others fit high-capacity networks to data without the use of such priors [28,32,29,26,25,19,31]. The latter methods additionally differ by the type of data structure used to represent the geometry, namely, mesh-based [9,26,25,12], point-based [28,51], and implicit models [29,32,35,39,31,27,50]. Additionally, recently there have emerged the hybrid models [55,10] where authors integrate face priors from parametric models with implicit representations to learn geometry and rendering for the specific person from the video.

However, mesh-based models arguably represent the most convenient class of methods for downstream applications. They provide better rendering quality and better temporal stability than point-based neural rendering. Also, unlike methods based on implicit geometry, mesh-based methods preserve topology and rigging capabilities and are much faster during fitting and rendering. However, current mesh-based methods either severely limit the range of deformations [9], making it infeasible to learn more complex geometry like hair and clothed shoulders, operate in the multi-shot setting [12] or require 3D scans as training data [26,25]. Our proposed method is also mesh-based, but we allow the prediction of complex deformations without 3D supervision and using a single image, lifting the limitations of the previous and concurrent works.

One-shot neural head models. Advances in neural networks also led to the development of methods that directly predict images using large ConvNets operating in the 2D image domain, with effectively no underlying 3D geometry [42,53,52] or very coarse 3D geometry [49]. These methods achieve state-of-the-art realism [49], use in-the-wild images or videos with no 3D annotations for training, and can create avatars from a single image. However, the lack of an explicit geometric model makes these models incompatible with many real-world applications and limits the span of camera poses that these methods can handle.

Neural mesh rendering. Recently, approaches that combine explicit data structures (point clouds or meshes) with neural image generation have emerged. These methods gained popularity thanks to the effectiveness of the pioneering Deferred Neural Rendering system [46], as well as recent advances in differen-

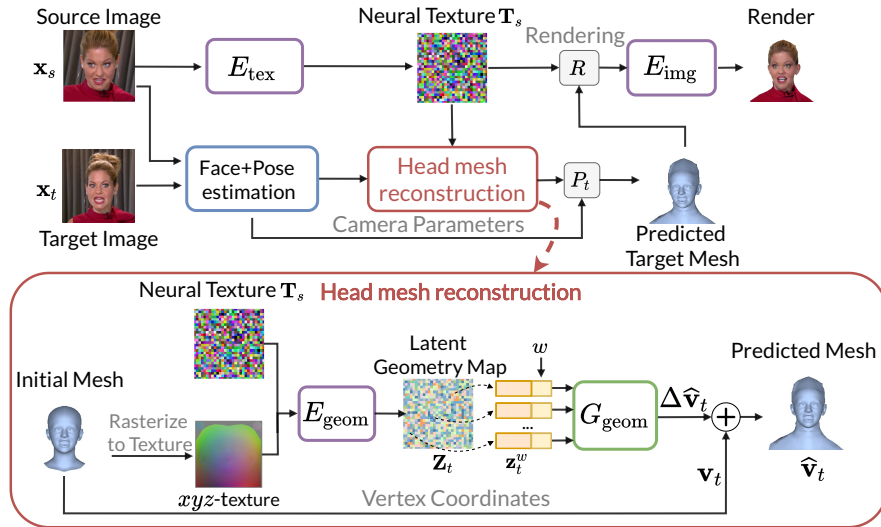


Fig. 2: Overview of our approach and the detailed scheme of the **head mesh reconstruction**. Given the source photo, we first estimate a *neural texture* that encodes local geometric and photometric details of visible and occluded parts. We then use a pre-trained system [9] for face reconstruction to estimate an initial mesh with a reconstructed facial part. We call this step face and 3D pose estimation. During head mesh reconstruction (bottom), using the estimated neural texture and the initial mesh, we predict the offsets for the mesh vertices, which do not correspond to a face. The offsets are predicted with a combination of a convolutional network E_{geom} and a perceptron network G_{geom} . We then render the personalized head mesh using the camera parameters, estimated by a pre-trained regressor [9] while superimposing the predicted neural texture. Finally, the rendering network E_{img} estimates the RGB image and the mask from the render.

table mesh rendering [36,24,22]. Neural mesh rendering uses 2D convolutional networks to model complex photometric properties of surfaces. It achieves high realism of renders with fine details present even when they are missing in the underlying geometric model. In this work, we adapt these advances to human head modelling while training using a large dataset of in-the-wild videos.

3 Method

Our goal is to build a system that jointly learns to produce photorealistic renders of human heads and estimate their 3D meshes using only a *single image* without any 3D supervision.

To achieve that, we use a large-scale dataset [6] of in-the-wild videos with talking speakers. All frames in each video are assumed to depict the same person

in the same environment (defined by lighting, hairstyle, and person’s clothing). At each training step, we sample two random frames \mathbf{x}_s and \mathbf{x}_t from a random training video. Our goal is to reconstruct and render the target image $\hat{\mathbf{x}}_t$ given a) the personal details and the face shape extracted from the source image \mathbf{x}_s , as well as b) the head pose, the facial expression, and the camera pose estimated from the target image \mathbf{x}_t . The final reconstruction loss is backpropagated and used to update the parameters of the model’s components.

After training, we can create a personalized head model by estimating all parameters from a single image. This model can then be *animated* using face tracking parameters extracted from any talking head sequence and rendered from a range of viewpoints similar to those present in the training dataset (Figure 1).

3.1 Model overview

In our model, we jointly train multiple neural networks that perform rendering and mesh reconstruction. The training pipeline proceeds as follows (Figure 2):

Neural texture estimation. The source image \mathbf{x}_s is encoded into a neural texture \mathbf{T}_s , which describes both person-specific appearance and geometry. The encoding is done by a convolutional neural network E_{tex} .

Face and 3D pose estimation. In parallel, we apply a pre-trained DECA system [9] for face reconstruction to both the source and the target image, which estimates facial shape, expression, and head pose. Internally, it uses the FLAME parametric head model [23], which includes mesh topology, texture mapping, and blendshapes. We use the face shape from the source image \mathbf{x}_s as well as the facial expression and the camera pose from the target image \mathbf{x}_t for further processing.

Head mesh reconstruction. The vertices of the DECA mesh with personalized face region and generic non-facial parts are rendered into an xyz -coordinate texture using the predefined texture mapping. The xyz -texture and the neural texture \mathbf{T}_s are concatenated and processed with the U-Net network [37] E_{geom} into a new texture map \mathbf{Z}_t , called *latent geometry* map. The 3D displacements for each mesh vertex are then decoded independently by the multi-layer perceptron G_{geom} that predicts a 3D offset $\Delta\hat{\mathbf{v}}$ for each vertex. This step reconstructs the personalized model for non-face parts of the head mesh. The obtained reconstructions are compatible with the topology/connectivity of the FLAME mesh [23].

Deferred neural rendering. The personalized head mesh is rendered using the pose estimated by DECA for the target image and with the superimposed neural texture. The rendered neural texture and the rasterized surface normals are concatenated and processed by the decoding (rendering) U-Net network E_{img} to predict the rendered image $\hat{\mathbf{x}}_t$ and the segmentation mask $\hat{\mathbf{s}}_t$. During training, the difference between the predictions and the ground truths is used to update all components of our system.

Below we discuss our system and its training process in more detail. We also describe a training procedure for a simplified version of our model, which represents head geometry using a linear basis of blendshapes.

3.2 Parametric face modeling

Our method uses a predefined head mesh with the corresponding topology, texture coordinates w , and rigging parameters, which remain fixed for all avatars. More specifically, we use FLAME [23] head model that has N base vertices $\mathbf{v}_{\text{base}} \in \mathbb{R}^{3N}$, and two sets of K and L basis vectors (blendshapes) that encode shape $\mathcal{B} \in \mathbb{R}^{3N \times K}$ and expression $\mathcal{D} \in \mathbb{R}^{3N \times L}$. The reconstruction process is carried out in two stages. First, the basis vectors are blended using the person- and expression-specific vectors of linear coefficients ϕ and ψ . Then, the linear blend skinning [23] function \mathcal{W} is applied, parameterized by the angles θ , which rotates the predefined groups of vertices around linearly estimated joints. The final reconstruction in world coordinates can be expressed as follows:

$$\mathbf{v}(\phi, \psi, \theta) = \mathcal{W}(\mathbf{v}_{\text{base}} + \mathcal{B}\phi + \mathcal{D}\psi, \theta).$$

In previous works [45], a similar set of parameters for the 3DMM [3] parametric model was obtained via photometric optimization. More recently, learning-based methods [9,13] capable of feed-forward estimation started to emerge. In our work, given an input image, we use a pre-trained feed-forward DECA system [9] to estimate ϕ, ψ, θ , and the camera parameters.

During training, we apply DECA to both source image \mathbf{x}_s and the target image \mathbf{x}_t . The face shape parameters ϕ_s from the source image \mathbf{x}_s alongside the expression ψ_t , head pose θ_t and camera parameters from the target image \mathbf{x}_t are then used to reconstruct the initial FLAME vertices $\mathbf{v}_t = \mathbf{v}(\phi_s, \psi_t, \theta_t)$, as well as camera transform \mathcal{P}_t .

3.3 Head mesh reconstruction

The FLAME vertices \mathbf{v}_t estimated by DECA provide good reconstructions for the face region but lack any person-specific details in the remaining parts of the head (hair and shoulders). To alleviate that, we predict person-specific mesh offsets for non-facial regions while preserving the face shape predicted by DECA. We additionally exclude ear regions since their geometry in the initial mesh is too complex to be learned from in-the-wild video datasets.

These mesh offsets are estimated in two steps. First, we encode both the xyz -coordinate texture and the neural texture \mathbf{T}_s into the latent geometry texture map \mathbf{Z}_t via a U-Net network E_{geom} . It allows the produced latent map to contain both positions of the initial vertices \mathbf{v}_t and their semantics, provided by the neural texture.

From \mathbf{Z}_t we obtain the vectors \mathbf{z}_t^w by bilinear interpolation at the fixed texture coordinates w . The vectors \mathbf{z}_t^w and their coordinates w are then concatenated and passed through a multi-layer perceptron G_{geom} to predict the coefficients $\hat{\mathbf{m}}_t \in \mathbb{R}^{3N \times 3}$ independently for each vertex in the mesh. These coefficients are multiplied elementwise by the normals \mathbf{n}_t , calculated for each vertex in \mathbf{v}_t , to obtain the displacements: $\Delta \hat{\mathbf{v}}_t = \hat{\mathbf{m}} \odot \mathbf{n}_t$. These displacements are then zeroed out for face and ear regions, and the final reconstruction in world coordinates is obtained as follows: $\hat{\mathbf{v}}_t = \mathbf{v}_t + \Delta \hat{\mathbf{v}}_t$.

3.4 Deferred neural rendering

We render the reconstructed head vertices $\hat{\mathbf{v}}_t$ using the topology and texture coordinates w from the FLAME model with the superimposed neural texture \mathbf{T}_s . For that, we use a differentiable mesh renderer \mathcal{R} [36] with the camera transform \mathcal{P}_t estimated by DECA for the target image \mathbf{x}_t .

The resulting rasterization, which includes both the neural texture and the surface normals, is processed by the rendering network E_{img} to obtain the predicted image $\hat{\mathbf{x}}_t$ and the segmentation mask $\hat{\mathbf{s}}_t$. E_{img} consists of two U-Nets that separately decode an image and a mask. The result of the deferred neural rendering is the reconstruction of the target image $\hat{\mathbf{x}}_t$ and its mask $\hat{\mathbf{s}}_t$, which are compared to the ground-truth image \mathbf{x}_t and mask \mathbf{s}_t respectively.

3.5 Training objectives

In our approach, we learn the geometry of hair and shoulders, which are not reconstructed by the pre-trained DECA estimator, without any ground-truth 3D supervision during training. For that we utilize two types of objectives: segmentation-based geometric losses $\mathcal{L}_{\text{geom}}$ and photometric losses $\mathcal{L}_{\text{photo}}$.

We found that explicitly assigning subsets of mesh vertices to the neck and the hair regions helps a lot with the quality of final deformations. It allows us to introduce a topological prior for the predicted offsets, which is enforced by.

To evaluate the geometric losses, we calculate two separate occupancy masks using a soft rasterization operation [24]. First, $\hat{\mathbf{o}}_t^{\text{hair}}$ is calculated with detached neck vertices, so that the gradient flows through that mask only to the offsets corresponding to the hair vertices, and then $\hat{\mathbf{o}}_t$ is calculated with detached hair vertices. We match the hair occupancy mask to the ground-truth mask $\mathbf{s}_t^{\text{hair}}$ (which covers the hair, face, and ears), and the estimated occupancy mask to the whole segmentation mask \mathbf{s}_t : $\mathcal{L}_{\text{occ}} = \lambda_{\text{hair}} \|\hat{\mathbf{o}}_t^{\text{hair}} - \mathbf{s}_t^{\text{hair}}\|_2^2 + \lambda_o \|\hat{\mathbf{o}}_t - \mathbf{s}_t\|_2^2$.

We also use an auxiliary Chamfer loss to ensure that the predicted mesh vertices cover the head more uniformly. Specifically, we match the 2D coordinates of the mesh vertices projected into the target image to the head segmentation mask. We denote the subset of predicted mesh vertices, visible in the target image, as $\hat{\mathbf{p}}_t = \mathcal{P}'_t(\hat{\mathbf{v}}_t)$, and the number of these vertices as N_t , so that $\hat{\mathbf{p}}_t \in \mathbf{R}^{N_t \times 2}$. Notice that operator \mathcal{P}'_t here not only does the camera transformation but also discards the z coordinate of the projected mesh vertices. To compute the loss, we then sample N_t 2D points from the segmentation mask \mathbf{s}_t and estimate the Chamfer distance between the sampled set of points \mathbf{p}_t and the vertex projections:

$$\mathcal{L}_{\text{chm}} = \frac{1}{2N_t} \sum_{\hat{p}_t \in \hat{\mathbf{P}}_t} \left\| \hat{p}_t - \arg \min_{p \in \mathbf{P}_t} \|p - \hat{p}_t\| \right\| + \frac{1}{2N_t} \sum_{p_t \in \mathbf{P}_t} \left\| p_t - \arg \min_{\hat{p} \in \hat{\mathbf{P}}_t} \|\hat{p} - p_t\| \right\|.$$

Lastly, we regularize the learned geometry using the Laplacian penalty [43]. Initially, we found that regularizing offsets $\Delta\hat{\mathbf{v}}$ worked better than regularizing full coordinates $\hat{\mathbf{v}}$ and stuck to that approach for all experiments. Our version of the Laplacian loss can be written as:

$$\mathcal{L}_{\text{lap}} = \frac{1}{V} \sum_{i=1}^V \left\| \Delta\hat{\mathbf{v}}_i - \frac{1}{|\mathcal{N}(i)|} \sum_{j \in \mathcal{N}(i)} \Delta\hat{\mathbf{v}}_j \right\|_1,$$

where $\mathcal{N}(i)$ denotes a set indices for vertices adjacent to the i -th vertex in the mesh.

We also use photometric optimization that matches the predicted and the ground truth images. This allows us to obtain photorealistic renders and aid in learning proper geometric reconstructions. We utilize perceptual loss \mathcal{L}_{per} [18], the face recognition loss \mathcal{L}_{idt} [5] and adversarial loss \mathcal{L}_{adv} [11,48]. We also use the Dice loss \mathcal{L}_{seg} [30] to match the predicted segmentation masks.

The final objective is weighted sum of the geometric and the photometric losses described above.

3.6 Linear deformation model

In addition to the full non-linear model introduced above, we consider a simplified parametric model with a linear basis of offsets $\Delta\hat{\mathbf{v}}$. While this model is similar to parametric models [23,56], we still do not use 3D scans for training and instead obtain our linear model by “distilling” the non-linear model. Additionally, we train a feed-forward estimator that predicts the linear coefficients from the input image.

The motivation for training this additional model is to show that the deformations learned by our method can be approximated using a system with a significantly lower capacity. Such a simple regression model can be easier to apply for inference on low-performance devices.

To train the linear model, we first obtain the basis of offsets $\mathcal{F} \in \mathbb{R}^{3N \times K}$, which is similar to the blendshapes used in the FLAME parametric model. This basis is obtained by applying a low-rank PCA [14] to the matrix of offsets $\Delta\hat{\mathbf{V}} \in \mathbb{R}^{3N \times M}$, calculated using M images from the dataset. Following [23], we discard most of the basis vectors and only keep K components corresponding to maximal singular values. The approximated vertex offsets $\tilde{\mathbf{v}}$ for each image can then be estimated as following $\tilde{\mathbf{v}} = \mathcal{F}\eta$, where η is obtained by applying the pseudo-inverse of a basis matrix \mathcal{F} to the corresponding offsets $\Delta\hat{\mathbf{v}}$: $\eta = (\mathcal{F}^T \mathcal{F})^{-1} \mathcal{F}^T \Delta\hat{\mathbf{v}}$

We then train the regression network by estimating a vector of basis coefficients η_t , given an image \mathbf{x}_t . For that, we minimize the mean squared error (MSE) loss $\|\hat{\eta}_t - \eta_t\|_2^2$ between the estimated coefficients and the ground truth, as well as the segmentation loss \mathcal{L}_{occ} and a Chamfer distance between predicted and ground truth meshes.

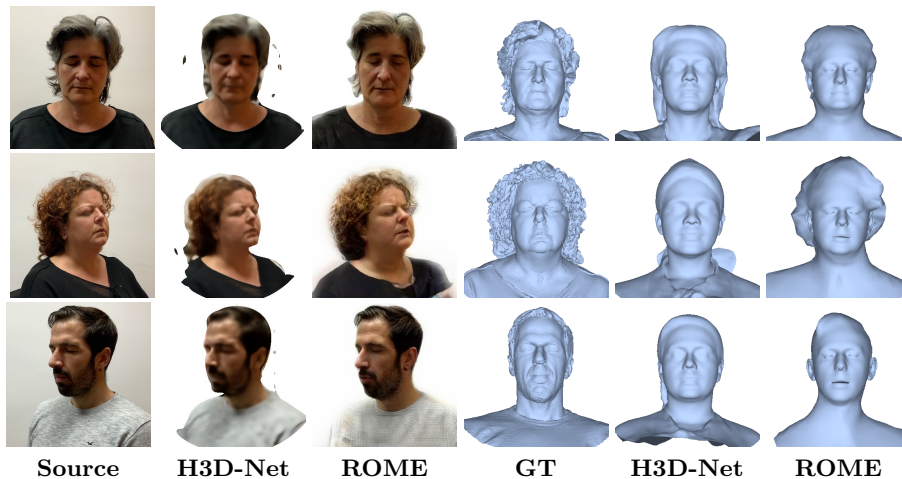


Fig. 3: Qualitative comparison of the representative cases from the H3DS dataset. While neither of the two methods achieves perfect results, arguably, ROME achieves more realistic renders and better matches the head geometry than H3D-Net in the single-shot mode. Furthermore, an important advantage of ROME is that the resulting avatars are ready for animation and are obtained in a feed-forward manner without the lengthy fine-tuning process employed by H3D-Net.

4 Experiments

We train our models on the VoxCeleb2 [6] dataset of videos. This large-scale dataset contains an order of 10^5 videos of 10^3 different speakers. It is widely used [7,49,52] to train talking head models. However, the main drawback of this dataset is the mixed quality of videos and the heavy bias towards frontal poses.

To address these well-known limitations, we process this dataset using an off-the-shelf image quality analysis model [44] and a 3D face-alignment network [4]. We then filter out the data which has poor quality and non-diverse head rotations. Our final training dataset has ≈ 15000 sequences. We note that filtering/pruning does not fully solve the problem of head rotation bias, and our method still works best in frontal views. For more details, please refer to the supplementary materials.

We also use the H3DS [35] dataset of photos with associated 3D scans to evaluate the quality of head reconstructions.

4.1 Implementation details

In the experiments, unless noted otherwise, we train all architectures jointly and end-to-end. We use the following weights: $\lambda_{\text{hair}} = 10$, $\lambda_{\text{per}} = 1$, $\lambda_{\text{idt}} = 0.1$, $\lambda_{\text{adv}} = 0.1$, $\lambda_{\text{seg}} = 10$, and enable the neck and the 2D Chamfer loss $\lambda_{\text{chm}} = 0.01$ and $\lambda_{\text{lap}} = 10$. We ablate all geometry losses and method parts below.

Table 1: Evaluation results on the H3DS dataset in the one-shot scenario for our models, H3D-Net, and DECA. We compute Chamfer distance (lower is better) across all available scans, reconstructed from three different viewpoints. Both of the ROME variants significantly exceed H3D-Net in the one-shot reconstruction quality.

Method	DECA	H3D-Net	ROME	LinearROME
Chamfer Distance	15.0	15.1	12.6	12.5

We train our models at 256×256 resolution using ADAM [21] with the fixed learning rate of 10^{-4} , $\beta_1 = 0$, $\beta_2 = 0.999$, and a batch size of 32. For more details, please refer to the supplementary materials.

4.2 Evaluation

3D reconstruction. We evaluate our head reconstruction quality using a novel H3DS dataset [35]. We compare against the state-of-the-art head reconstruction method H3D-Net [35], which uses signed distance functions to represent the geometry. While providing great reconstruction quality in the sparse-view scenario, their approach has several limitations. For example, H3D-Net requires a dataset of full head scans to learn the prior on head shapes. Additionally, its results do not have fixed topology or rigging and their method requires fine-tuning per scene, while our method works in a feed-forward way.

We carry out the comparison with H3D-Net in a single-view scenario, which is native for our method but is beyond the capabilities stated by the authors in the original publication [35]. However, to the best of our knowledge, H3D-Net is the closest system to ours in single-view reconstruction capabilities (out of all systems with either their code or results available). Additionally, we tried to compare our system with PIFuHD [38], which unfortunately failed to work with heads images without body (see supplementary).

We show qualitative comparison in Figure 3. We evaluate our method and H3D-Net both for frontal- and side-view reconstruction. We note the significant overfitting of H3D-Net to the visible hair geometry, while our model provides reconstructions more robust to the change of viewpoint.

In total, we compared our models on all scans available in the test set of the H3DS dataset, and each scan was reconstructed from three different viewpoints. We provide the measured mean Chamfer distance both for our method and baselines across all scans in Tab. 1.

Rendering. We evaluate the quality of our renders on the hold-out subset Vox-Celeb2 dataset. We use a cross-driving comparison scenario for qualitative comparison to highlight the animation capabilities of our method, and self-driving scenario for quantitative comparison.

First, we compare with a FLAMETex [23] rendering system, which works explicitly with mesh rendering. From the source image, FLAMETex estimates the albedo via a basis of RGB textures, and then combines it with predicted



Fig. 4: Comparison of renders on a VoxCeleb2 dataset. The task is to reenact the **source** image with the expression and pose of the **driver** image. Here, we picked diverse examples in terms of pose variation to highlight the differences in performance of compared methods. We observe that for the large head pose rotations, purely neural-based methods (**FOMM**, **Bi-Layer**) struggle to maintain consistent quality. In contrast, our rendering method produces images that are more robust to pose changes. Admittedly, for small pose changes, neural-based methods exhibit a smaller identity gap than ROME (bottom row) and overall outperform our method in terms of rendering quality. As a reference, we also include a non-neural **FLAMETex** rendering method, which is employed in state-of-the-art one-shot face reconstruction systems [9] but is not able to personalize the avatar at the head level.

scene-specific shading. On the contrary, our method predicts a rendered image directly and avoids the complexity of explicit albedo-shading decomposition.

We then compare with publicly available geometry-free rendering methods, which were trained on the same dataset. For that, we use the First-Order Motion Model (FOMM) [42], the Bi-Layer Avatar Model [52] and recently proposed Thin-Plate-Spline-Motion-Mode (TPSMM) [54]. Both these systems bypass explicit 3D geometry estimation and rely only on learning the scene structure via the parameters of generative ConvNets. Other methods [49,7], which internally utilize some 3D structures, like camera rotations, were out of the scope of our comparison due to the unavailability of pre-trained models.

We present the qualitative comparison in Figure 4, and a quantitative comparison across a randomly sampled hold-out VoxCeleb2 subset in Table 2. We restrict the comparison to the face and hair region as the shoulder pose is not controlled by our method (driven by DECA parameters), which is admittedly a limitation of our system. We thus mask the results according to the face and hair mask estimated from the ground truth image.

Table 2: Here we present the quantitative results on the VoxCeleb2-HQ dataset in the self-reenactment and cross-reenactment modes. Our ROME system performs on par with FOMM and TPSMM in self-reenactment, notably outperforming them in the most perceptually-plausible LPIPS metrics. On the contrary, in the cross-driving scenario, when the task is complex for pure neural-based systems, our method obtains better results.

Method	self-reenactment			cross-reenactment		
	LPIPS↓	SSIM↑	PSNR↑	FID↓	CSIM↑	IQA↑
FOMM	0.09	0.87	25.8	52.95	0.53	55.9
Bi-Layer	0.08	0.83	23.7	51.4	0.56	50.48
TPSMM	0.09	0.85	26.1	49.27	0.57	59.5
ROME	0.08	0.86	26.2	45.32	0.62	66.3

Generally, we observe that over the entire test set, the quality of ROME avatars in the self-reenactment mode is similar to FOMM and better than the Bi-layer model. For the cross-reenactment scenario, our model is clearly better than both alternatives according to three metrics, that help to assess unsupervised quality of the images in three aspects: realism, identity preservation and blind quality of the image. The huge gap for IQA [44] and FID [17] is also noticeable in the qualitative comparison, especially for strong pose change (see CSIM [53] column in Tab. 2). The PSNR and SSIM metrics penalize slight misalignments between the sharp ground truth and our renderings much stronger than the blurriness in FOMM reconstructions. The advantage of ROME avatar is noticeable even for self-driving case according to LPIPS. We provide a more extensive qualitative evaluation in the supplementary materials.

4.3 Linear basis experiments

As discussed above, we distill our ROME head reconstruction model into a linear parametric model. To do that, we set the number of basis vectors to 50 for the hair and 10 for the neck offsets and run low-rank Principle Component Analysis (PCA) to estimate them. The number of components is chosen to obtain a low enough approximation error. Interestingly, the offsets learned by our model can be compressed by almost two orders of magnitude in terms of degrees of freedom without any practical loss in quality (Figure 6a), which suggests that the capacity of the offset generator is underused in our model. We combine estimated basis vectors with the original basis of the FLAME.

After that, we train feed-forward encoders that directly predict the coefficients of the two basis from the source image. The prediction is performed in two stages. First, face expression, pose and camera parameters are predicted with a MobileNetV2 [40] encoder. Then a slower ResNet-50 encoder [16] is used to predict hair, neck and shape coefficients. The choice of architectures are motivated by the fact that in many practical scenarios only the first encoder needs

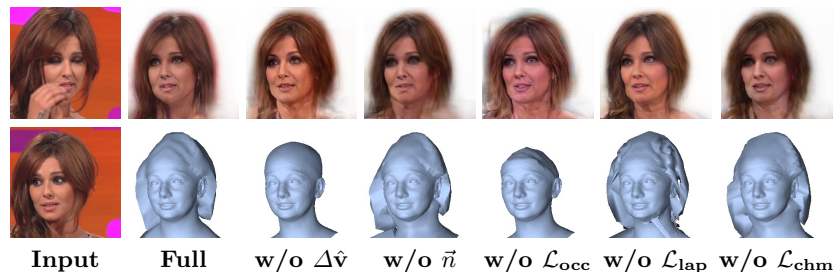


Fig. 5: Ablation study. We qualitatively evaluate the individual components of our *full* model. *w/o* $\Delta\hat{\mathbf{v}}$: without the per-vertex displacements, we obtain a significantly worse render quality. *w/o* $\bar{\mathbf{n}}$: when we apply per-vertex deformations instead of per-vertex displacements (i.e., deformations alongside the normals), we obtain noisy reconstructions in neck area and worse renders. *w/o* \mathcal{L}_{occ} : without silhouette-based losses, our model fails to learn proper reconstructions. *w/o* \mathcal{L}_{lap} : Laplacian regularization smooths the reconstructions. *w/o* \mathcal{L}_{chm} : chamfer loss allows us to constrain the displaced vertices to lie inside the scene boundaries, which positively affects the smoothness of the visible part of the reconstruction.

to be invoked frequently (per-frame), while the second can run at much lower rate or even only at the model creation time.

4.4 Ablation study

We demonstrate results of ablation study at Figure 5. As expected, predicting more accurate geometry affect the renders (first row). Also, we verify the necessity of all terms of geometry loss. We observe significant improvement in quality of renders with additional geometry (see supplementary), which leads us to an optimistic conclusion that our learned coarse mesh may be integrated into other neural rendering systems [49] to improve their quality. Additionally, we observe that geometry losses allows to correctly model coarse details on the hair without noise and reconstruct the hair without sticking with neck. Similar artifacts are removed by adding shifts along the normals.

Our current model is trained at roughly fixed scale, though explicit geometry modeling allows it to generalize to adjacent scale reasonably well. Still, strong changes of scale lead to poor performance (Figure 6b). More examples are provided in the supplementary materials. Addressing this issue via mip-mapping and multi-scale GAN training techniques remains future work.

Lastly, our model can have artifacts with long hair (Figure 6b, left) or ears (Figure 6b, middle). Handling such cases gracefully are likely to require a departure from the predefined FLAME mesh connectivity to new person-specific mesh topology. Handling such issues using a limited set of pre-designed hair meshes is an interesting direction for future research.

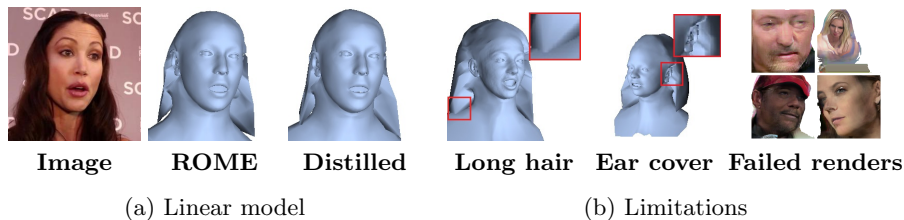


Fig. 6: Linear model results and the examples of limitations. On the left, we show how reconstructions learned by our method, **ROME**, could be **distilled** using a linear parametric model. We are able to compress the learned offsets into a small basis, reducing the degrees of freedom by two orders of magnitude. We can then **distill** these offsets using a much faster regression network with a small gap in terms of the reconstruction quality. On the right, we highlight the main limitations of our method, which include the failure related to **long hair** modelling, caused by an incorrect topological prior, no **coverage of ears** and **unrealistic renders** under a significant change of scales.

5 Summary

We have presented ROME avatars: a system for creating realistic one-shot mesh-based human head models that can be animated and compatible with FLAME head models. We compare our model with representative state-of-the-art models from different classes, and show that it is highly competitive both in terms of geometry estimation and the quality of rendering.

Crucially, our system can learn to model head geometry without direct supervision in the form of 3D scans. Despite that, we have observed it to achieve state-of-the-art results in head geometry recovery from a single photograph. At the same time, it also performs better than previous one-shot neural rendering approaches in the cross- and self-driving scenario. We have thus verified that the resulting geometry could be used to improve the rendering quality.

As neural rendering becomes more widespread within graphics systems, ROME avatars and similar systems can become directly applicable, while their one-shot capability and the simplicity of rigging derived from DECA and FLAME could become especially important in practical applications.

Acknowledgements

We sincerely thank Eduard Ramon for providing us the one-shot H3D-Net reconstructions. We also thank Arsenii Ashukha for comments and suggestions regarding the text contents and clarity, as well as Julia Churkina for helping us with proof-reading. The computational resources for this work were mainly provided by Samsung ML Platform.

References

1. AvatarSDK: <https://avatarsdk.com/> **2**
2. Blanz, V., Romdhani, S., Vetter, T.: Face identification across different poses and illuminations with a 3d morphable model. *Proceedings of Fifth IEEE International Conference on Automatic Face Gesture Recognition* pp. 202–207 (2002) **2**
3. Blanz, V., Vetter, T.: A morphable model for the synthesis of 3d faces. In: *SIGGRAPH '99* (1999) **1, 2, 6**
4. Bulat, A., Tzimiropoulos, G.: How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In: *International Conference on Computer Vision* (2017) **9**
5. Cao, Q., Shen, L., Xie, W., Parkhi, O.M., Zisserman, A.: Vggface2: A dataset for recognising faces across pose and age. *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)* pp. 67–74 (2018) **8**
6. Chung, J.S., Nagrani, A., Zisserman, A.: Voxceleb2: Deep speaker recognition. In: *INTERSPEECH* (2018) **2, 3, 4, 9**
7. Doukas, M.C., Zafeiriou, S., Sharmanska, V.: Headgan: Video-and-audio-driven talking head synthesis (2021) **9, 11**
8. Egger, B., Smith, W., Tewari, A., Wuhler, S., Zollhöfer, M., Beeler, T., Bernard, F., Bolkart, T., Kortylewski, A., Romdhani, S., Theobalt, C., Blanz, V., Vetter, T.: 3d morphable face models—past, present, and future. *ACM Transactions on Graphics (TOG)* **39**, 1 – 38 (2020) **2**
9. Feng, Y., Feng, H., Black, M.J., Bolkart, T.: Learning an animatable detailed 3d face model from in-the-wild images. *ACM Transactions on Graphics (TOG)* **40**, 1 – 13 (2020) **1, 2, 3, 4, 5, 6, 11**
10. Gafni, G., Thies, J., Zollhöfer, M., Nießner, M.: Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2021) **3**
11. Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A.C., Bengio, Y.: Generative adversarial nets. In: *NIPS* (2014) **8**
12. Grassal, P.W., Prinzel, M., Leistner, T., Rother, C., Nießner, M., Thies, J.: Neural head avatars from monocular rgb videos. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2022) **3**
13. Guo, J., Zhu, X., Yang, Y., Yang, F., Lei, Z., Li, S.Z.: Towards fast, accurate and stable 3d dense face alignment. In: *Proceedings of the European Conference on Computer Vision (ECCV)* (2020) **1, 2, 6**
14. Halko, N., Martinsson, P.G., Tropp, J.A.: Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Rev.* **53**, 217–288 (2011) **8**
15. Hassner, T., Harel, S., Paz, E., Enbar, R.: Effective face frontalization in unconstrained images. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* pp. 4295–4304 (2015) **2**
16. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* pp. 770–778 (2016) **12**
17. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. In: *Advances in Neural Information Processing Systems* (2017) **12**
18. Johnson, J., Alahi, A., Fei-Fei, L.: Perceptual losses for real-time style transfer and super-resolution. In: *ECCV* (2016) **8**

19. Kellnhöfer, P., Jebe, L., Jones, A., Spicer, R.P., Pulli, K., Wetzstein, G.: Neural lumigraph rendering. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2021) [3](#)
20. Kim, H., Garrido, P., Tewari, A., Xu, W., Thies, J., Nießner, M., Pérez, P., Richardt, C., Zollöfer, M., Theobalt, C.: Deep video portraits. *ACM Transactions on Graphics (TOG)* **37**(4), 163 (2018) [2](#)
21. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *International Conference for Learning Representations* (2015) [10](#)
22. Laine, S., Hellsten, J., Karras, T., Seol, Y., Lehtinen, J., Aila, T.: Modular primitives for high-performance differentiable rendering. *ACM Transactions on Graphics* **39**(6) (2020) [4](#)
23. Li, T., Bolkart, T., Black, M.J., Li, H., Romero, J.: Learning a model of facial shape and expression from 4d scans. *ACM Transactions on Graphics (TOG)* **36**, 1 – 17 (2017) [2](#), [5](#), [6](#), [8](#), [10](#)
24. Liu, S., Li, T., Chen, W., Li, H.: Soft rasterizer: A differentiable renderer for image-based 3d reasoning. 2019 IEEE/CVF International Conference on Computer Vision (ICCV) pp. 7707–7716 (2019) [4](#), [7](#)
25. Lombardi, S., Saragih, J.M., Simon, T., Sheikh, Y.: Deep appearance models for face rendering. *ACM Transactions on Graphics (TOG)* **37**, 1 – 13 (2018) [2](#), [3](#)
26. Lombardi, S., Simon, T., Saragih, J.M., Schwartz, G., Lehrmann, A.M., Sheikh, Y.: Neural volumes. *ACM Transactions on Graphics (TOG)* **38**, 1 – 14 (2019) [2](#), [3](#)
27. Lombardi, S., Simon, T., Schwartz, G., Zollhoefer, M., Sheikh, Y., Saragih, J.M.: Mixture of volumetric primitives for efficient neural rendering. *ACM Transactions on Graphics (TOG)* **40**, 1 – 13 (2021) [3](#)
28. Ma, Q., Saito, S., Yang, J., Tang, S., Black, M.J.: Scale: Modeling clothed humans with a surface codec of articulated local elements. In: *CVPR* (2021) [3](#)
29. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: *ECCV* (2020) [2](#), [3](#)
30. Milletari, F., Navab, N., Ahmadi, S.A.: V-net: Fully convolutional neural networks for volumetric medical image segmentation. 2016 Fourth International Conference on 3D Vision (3DV) pp. 565–571 (2016) [8](#)
31. Oechsle, M., Peng, S., Geiger, A.: Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. *International Conference on Computer Vision (ICCV)* (2021) [3](#)
32. Park, J.J., Florence, P.R., Straub, J., Newcombe, R.A., Lovegrove, S.: Deepsdf: Learning continuous signed distance functions for shape representation. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 165–174 (2019) [3](#)
33. Pinscreen: <https://www.pinscreen.com/> [2](#)
34. Ploumpis, S., Ververas, E., Sullivan, E.O., Moschoglou, S., Wang, H., Pears, N.E., Smith, W., Gecer, B., Zafeiriou, S.: Towards a complete 3d morphable model of the human head. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021) [2](#)
35. Ramon, E., Triginer, G., Escur, J., Pumarola, A., Giraldez, J.G., i Nieto, X.G., Moreno-Noguer, F.: H3d-net: Few-shot high-fidelity 3d head reconstruction. *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021) [2](#), [3](#), [9](#), [10](#)
36. Ravi, N., Reizenstein, J., Novotny, D., Gordon, T., Lo, W.Y., Johnson, J., Gkioxari, G.: Accelerating 3d deep learning with pytorch3d. *arXiv:2007.08501* (2020) [2](#), [4](#), [7](#)

37. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: MICCAI (2015) [5](#)
38. Saito, S., Simon, T., Saragih, J., Joo, H.: Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (June 2020) [10](#)
39. Saito, S., Simon, T., Saragih, J.M., Joo, H.: Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 81–90 (2020) [3](#)
40. Sandler, M., Howard, A.G., Zhu, M., Zhmoginov, A., Chen, L.C.: Mobilenetv2: Inverted residuals and linear bottlenecks. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (2018) [12](#)
41. Sanyal, S., Bolkart, T., Feng, H., Black, M.J.: Learning to regress 3d face shape and expression from an image without 3d supervision. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 7755–7764 (2019) [1](#)
42. Siarohin, A., Lathuilière, S., Tulyakov, S., Ricci, E., Sebe, N.: First order motion model for image animation. Advances in Neural Information Processing Systems (NeurIPS) (2019) [2](#), [3](#), [11](#)
43. Sorkine-Hornung, O.: Laplacian mesh processing. In: Eurographics (2005) [8](#)
44. Su, S., Yan, Q., Zhu, Y., cui Zhang, C., Ge, X., Sun, J., Zhang, Y.: Blindly assess image quality in the wild guided by a self-adaptive hyper network. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2020) [9](#), [12](#)
45. Thies, J., Zollhöfer, M., Stamminger, M., Theobalt, C., Nießner, M.: Face2face: Real-time face capture and reenactment of rgb videos. In: Proc. Computer Vision and Pattern Recognition (CVPR), IEEE (2016) [2](#), [6](#)
46. Thies, J., Zollhöfer, M., Nießner, M.: Deferred neural rendering: Image synthesis using neural textures. arXiv: Computer Vision and Pattern Recognition (2019) [2](#), [3](#)
47. Tran, A., Hassner, T., Masi, I., Medioni, G.G.: Regressing robust and discriminative 3d morphable models with a very deep neural network. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) pp. 1493–1502 (2017) [2](#)
48. Wang, T.C., Liu, M.Y., Zhu, J.Y., Tao, A., Kautz, J., Catanzaro, B.: High-resolution image synthesis and semantic manipulation with conditional gans. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition pp. 8798–8807 (2018) [8](#)
49. Wang, T.C., Mallya, A., Liu, M.Y.: One-shot free-view neural talking-head synthesis for video conferencing. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 10034–10044 (2021) [2](#), [3](#), [9](#), [11](#), [13](#)
50. Yenamandra, T., Tewari, A., Bernard, F., Seidel, H.P., Elgharib, M.A., Cremers, D., Theobalt, C.: i3dmm: Deep implicit 3d morphable model of human heads. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 12798–12808 (2021) [3](#)
51. Zakharkin, I., Mazur, K., Grigoriev, A., Lempitsky, V.S.: Point-based modeling of human clothing. 2021 IEEE/CVF International Conference on Computer Vision (ICCV) (2021) [3](#)
52. Zakharov, E., Ivakhnenko, A., Shysheya, A., Lempitsky, V.S.: Fast bi-layer neural synthesis of one-shot realistic head avatars. In: ECCV (2020) [2](#), [3](#), [9](#), [11](#)
53. Zakharov, E., Shysheya, A., Burkov, E., Lempitsky, V.S.: Few-shot adversarial learning of realistic neural talking head models. In: 2019 IEEE/CVF International Conference on Computer Vision (ICCV) (2019) [2](#), [3](#), [12](#)

54. Zhao, J., Zhang, H.: Thin-plate spline motion model for image animation. 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2022) 11
55. Zheng, Y., Abrevaya, V.F., Bühler, M.C., Chen, X., Black, M.J., Hilliges, O.: I m avatar: Implicit morphable head avatars from videos. In: 2022 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2022) 3
56. Zuffi, S., Kanazawa, A., Berger-Wolf, T., Black, M.J.: Three-d safari: Learning to estimate zebra pose, shape, and texture from images ”in the wild” (2019) 8