

Leveraging Deepfakes to Close the Domain Gap between Real and Synthetic Images in Facial Capture Pipelines

Winnie Lin¹, Yilin Zhu^{1,2}, Demi Guo^{1,2}, and Ron Fedkiw^{1,2}

¹Stanford University

²Epic Games

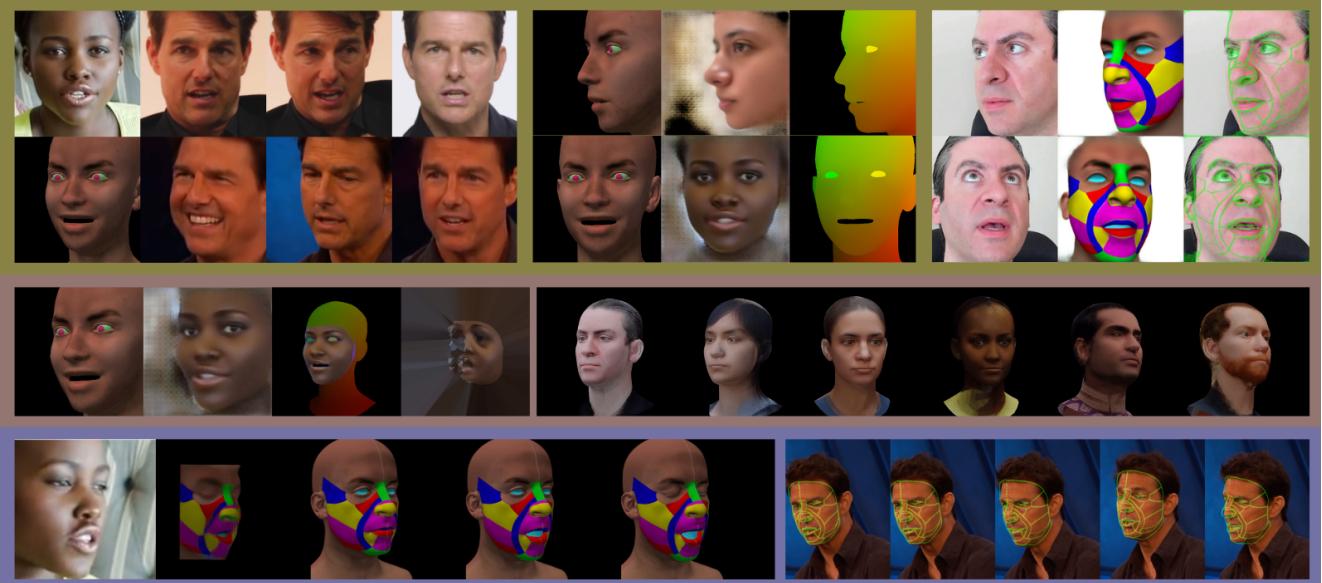


Figure 1. Top left: Given either real or synthetic input images, our automatic data curation allows us to find similar pose and expressions from an in-the-wild image dataset. This allows us to train personalized networks in a scalable way, requiring only a few hundred to a few thousand in-the-wild images collected via cellphones, webcams, or youtube videos. Top middle: The ability to inference from synthetic to real is used in our appearance capture pipeline. Top right: The ability to inference from real to synthetic is used in our motion tracking pipeline. Middle row: Our appearance capture pipeline, along with some representative results (on various races). Third row: Our motion capture pipeline, illustrating results from both inverse rendering and regression.

Abstract

We propose an end-to-end pipeline for both building and tracking 3D facial models from personalized in-the-wild (cellphone, webcam, youtube clips, etc.) video data. First, we present a method for automatic data curation and retrieval based on a hierarchical clustering framework typical of collision detection algorithms in traditional computer graphics pipelines. Subsequently, we utilize synthetic turntables and leverage deepfake technology in order to build a synthetic multi-view stereo pipeline for appearance capture that is robust to imperfect synthetic geometry and image misalignment. The resulting model is fit with an animation rig, which is then used to track facial performances. Notably, our novel use of deepfake technology enables us to perform

robust tracking of in-the-wild data using differentiable renderers despite a significant synthetic-to-real domain gap. Finally, we outline how we train a motion capture regressor, leveraging the aforementioned techniques to avoid the need for real-world ground truth data and/or a high-end calibrated camera capture setup.

1. Introduction

There is a growing interest in avatar personalization and facial motion tracking for consumer use, as avatars that accurately capture the likeness, appearance, and motion of the user increase user engagement and immersiveness. Accordingly, there is a plethora of industry interest in 3D avatar customization as part of a greater push for immersive enter-

tainment and community building, as demonstrated by the popularity of Roblox, Unreal Engine’s Metahuman creator, and various virtual/augmented reality related avatar personalization endeavours from Meta, Apple, and others. Even when a personalized avatar is not desired, they are quite useful as strong priors for the motion capture required for puppeteering.

Highly-specialized high-end methods such as lightstage capture [18], multi-view stereo reconstruction [4], head-mounted cameras [5] etc. have an important role in the special effects industry, enabling the acquisition of appearance, geometry, and motion from image (and marker) data, but these pipelines are computationally and resource intensive, and require carefully calibrated hardware, domain expertise, and manual intervention to operate. Meanwhile, with the advent of data-driven deep-learning methods, there is an exciting surge in the democratization of avatar creation and generation. Although this is still an actively developing area of research with many open-ended questions, it is undeniable that building semi-automated pipelines that do not require specialized setup and expertise is a key step to creating personalized avatars at scale.

Fully data-driven deep-learning methods often suffer from either single-view based input constraints, data diversity problems, or both. Single-view based appearance and motion capture is a wildly underconstrained problem, and it is still unclear how one would best incorporate arbitrary multi-view in-the-wild imagery into a pure deep-learning approach. The overreliance on large datasets gives rise to many practical problems regarding data quality, diversity, and bias. Although this is a common issue for many deep-learning applications, data diversity is particularly crucial when personalization/customization is of key importance. For example, papers in the AI ethics community [11] have pointed out that facial models trained on Caucasian-centric (or even worse, Caucasian male-centric) data perform poorly and unreliably on other ethnic groups. Problems like this are inherent in most state-of-the-art methods, and need to be adequately addressed [29][43] before such methods are able to be deployed at scale. To mitigate these issues, our approach overfits small, specialized deep-learning models to each specific subject of interest. In order to make such a method tractable and scalable enough for democratization, we automate the curation and compactification of personalized input data, relying on our proposed data curation algorithms to aggressively prune the data and thus reduce the compute time needed to train specialized deep neural networks.

Our pipeline aims to capture the strengths of both traditional methods and deep-learning methods. More specifically, for appearance capture, we rely on classical multi-view stereo techniques to be less beholden to the monocular-view constraints of deep-learning approaches, while utilizing synthetic turntables and leveraging deepfake technology to mit-

igate multi-view stereo misalignment and obviate the need for any specialized hardware. For motion capture, we utilize deepfake technology to close the domain gap between real in-the-wild images and synthetically-generated data, facilitating robust inverse-rendering even in the absence of photorealistic appearance and lighting models. In addition, a similar leveraging of deepfake technology allows us to train a fast and efficient motion capture regressor in the absence of ground truth control parameters for real imagery.

2. Prior Work

Existing approaches for appearance and motion capture typically fall into two categories: either high-end personalized approaches widely used in the visual effects industry or democratized monolithic approaches widely used in commercial settings. Here, we briefly summarize some relevant work, referring the interested reader to [64] for a more in-depth survey.

High-end methods: In [2], the authors built a multi-view stereo pipeline carefully optimized for geometry capture, reconstructing pore-level detail through iterative stereo-pair based geometric reconstruction. In [18], the authors built a controlled spherical lighting setup for high quality capture of the reflectance model of a face. Both of these methods and their successors (e.g. [23][52]) are an integral component of appearance capture pipelines in most major visual effects studios [17][25][3].

Early works on high-end motion capture reconstructed dense mesh sequences from video, with methods such as [63] and [42] employing multi-view stereo and scanline projectors. Current high-end motion capture approaches typically require both head-mounted cameras and physical markers placed on the face (see e.g. [16][45][9]). Solving the capture problem with sparse markers necessitates either a personalized animation rig (created via multi-view stereo [4], deformation transfer [55], artist supervision, etc.) or a data-driven approach (e.g. [16] interpolates from reconstructed and/or simulated geometry, [25] performs dense marker driven mesh deformations, etc.). Markerless methods, previously more commonly utilized in democratized approaches, have in recent years also become an active area of research for high-end applications, with inverse rendering technology (e.g. [51][38]) being used to solve for animation parameters directly from image data (e.g. [1][24]).

Democratized methods: These methods typically start with a parametrized template geometry/texture usable across all subjects. [7] pioneered a statistical PCA approach, generating a linear parametrization from a dataset of scans. This approach has been the standard for many years, with gradually increasing dataset size and complexity [8][34]. More recently, deep-learning approaches (which generate nonlinear parametrizations) have gained popularity (e.g. [56][13]), and hybrid PCA/deep-learning methods [62] have also emerged.

Both the PCA and deep-learning approaches typically parameterize both appearance and expression [20]; as such, they can be utilized for both appearance and motion capture.

To perform appearance capture, it is common to utilize facial landmarks (e.g. [10]) or other sparse features as loss constraints for training regressors that predict geometry parameters from monocular images (e.g. [53][57]). Democratized motion capture approaches typically utilize low-dimensional parametrizations of geometry, either data-driven (e.g. [7][56]) or artist-sculpted (e.g. [32]). The parameters are typically determined using optimization to fit the geometry directly to landmarks [14] or depth maps [61][33]. Alternatively, monolithic deep-learning regressors can be trained via sparse landmark constraints to directly output parameters from images [36][58][47]. A recent example of a regressor based approach [53] jointly predicts appearance and pose/expression in their end-to-end model, utilizing a landmark loss and a mesh based triplet loss during training.

Hybrid approaches: The development of commercially available AR/VR headsets has led to research on democratized markerless motion capture for head-mounted cameras, see e.g. [40][59]. These methods typically use monolithic multi-identity neural networks, although [28] adds a personalized conditioner by inputting a neutral face mesh of the subject. Our approaches to both appearance and motion capture similarly combine high-end and democratized approaches. We aim to build personalized pipelines typical of high-end applications while also dealing with the uncertainties of in-the-wild cameras and data.

Notably, [44] (contemporaneous with our work) takes a similar approach to training personalized motion capture regressors, learning a joint embedding space between real and synthetic images for markerless tracking in a seated motion capture setup. Although their approach does not handle in-the-wild data, it does show promising generalizability to less-constrained data capture. In addition, our work on the motion capture regressor is heavily related to and preceded by [37], which utilizes joint embeddings between real and synthetic images to do markerless motion capture with a multi-view head-mounted camera; however, their controlled capture setup is more stringent than what is required of our approach (and [44]). Our use of inverse rendering for offline motion capture also bears similarity to [22] and [58]. [22] describes a single-image based appearance capture pipeline, and utilizes a neural renderer for offline optimization. [58] uses a neural renderer to learn unconstrained deformable motion from in-the-wild images. In comparison, our inverse-rendering approach differentiates through a dense pixel loss instead of sparse keypoints and is a single-identity model trained only on subject-specific data, but unlike [58] our method is currently not designed to run in real-time.

Our approach to appearance capture is similar in spirit to [26][39] [46]. [26] utilizes a traditional model fitting

approach that generates avatars from a short video clip of a single subject; in our approach, we obviate the need for camera estimation and image alignment by utilizing deepfake networks. [39] achieves high-quality results on geometry/texture acquisition via a multi-identity GAN, similarly using dramatically less data than other contemporary works through careful selection and design of their dataset. [46] builds a multi-identity pipeline for face frontalization, where an image of a frontal face with diffuse lighting is generated from inputs consisting of arbitrary pose/expression and lighting. While leveraging similar concepts to [39][46], our method focuses on using a collection of images of a single subject instead of single images of multiple subjects.

3. Hierarchical Clustering for Data curation

A crucial part of both our appearance capture and motion capture pipelines is the curation of data from in-the-wild footage. The importance of data curation in our approach is twofold: We aim to not only prune irrelevant and noise-inducing data, but also to craft compact datasets with the exact amount of information required for a given application. As shown in section 5, our data curation allows us to successfully train deepfake networks from scratch with fewer than 200 images in each dataset. In order to curate compact datasets, we devised a k-d tree based [30] recursive partitioning of semantically meaningful image clusters (which additionally improves the efficiency of image retrieval.) Through iterative experimentation, we designed several features corresponding to pose, expression, and lighting:

Pose. The 3-dimensional pose of a rigid object such as the skull is most commonly expressed in terms of a three degrees of freedom for translation and three degrees of freedom for rotation. However, when considering the projection of a rigid object into a 2D image, one can make some rough approximations in order to remove some degrees of freedom. Assuming minimal lens distortion (depth distortion, field of view dependencies, etc.) allows one to remove all three degrees of freedom corresponding to translation via image cropping and scaling. Similar assumptions allow one to remove rotation along the camera look-at axis, instead orienting the face to be upright via image rotation. This leaves two pronounced degrees of freedom, pitch and yaw, which both create significant occlusion of facial features.

Given an image I of a face, we extract the pose feature $f_{pose}(I)$ as follows: we detect 68 2D facial landmarks $L(I) \in \mathbb{R}^{68 \times 2}$ [10] and identify a subset of these landmarks $L_r \in \mathbb{R}^{13 \times 2}$ that don't vary too much with expression. We rigidly align L_r to a predefined set of corresponding template landmarks \tilde{L}_r that were derived assuming a frontal view with zero pitch and yaw. Even though L_r mostly

comes from images with nonzero pitch and yaw, this rigid alignment helps reduce dependencies on both scale and in-plane affine transformations enough to fit a 3D model to the result. That is, we can then determine the pitch and yaw parameters of a 3D template model that minimize the distance between the rigidly aligned L_r and the 2D projection of corresponding markers embedded on the 3D model. Our pose feature $f_{pose}(I)$ is then the 2-dimensional vector corresponding to the resulting pitch and yaw of the 3-dimensional model. See Figure 2. The least squares problem for finding the optimal rigid transformations for both steps can be solved using a simple singular value decomposition on a 3×3 matrix [60][54].



Figure 2. From left to right: the image I , the image rigidly aligned to the set of template landmarks, a template 3D model rendered with the resulting pitch and yaw.

Expression. We take landmarks $L_e \subset L(I)$ that correspond more heavily to expression (rather than pose), split them into (potentially overlapping) groups, and align each group separately to a 2D template \tilde{L}_e in order to normalize for scale, in-plane rotation, and translation in a parts based manner. (Parts based models are common in human pose identification [12], and bag-of-words approaches have also been utilized for facial recognition [35].) The results are then directly used as our feature vectors. See Figure 3. Of course, this high-dimensional feature vector could benefit from k-d tree based acceleration structures, but in practice our expression clusters are small enough that we leave their storage flat.

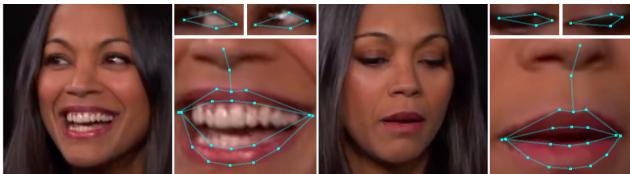


Figure 3. Image I and the aligned markers $f_{expression}(I)$ used in the expression matching layer (two examples are shown). To account for differences in global facial structure between a query subject and a dataset subject, we use a part based approach aligning the individual eyes and mouth separately.

An important technical detail is that the efficacy of $f_{expression}(I)$ hinges on an initial match based on pose, i.e. poor matches are common when the rigid poses are dissimilar. We address this by ensuring that $f_{expression}(I)$

only comes after $f_{pose}(I)$ in our hierarchical approach.

Lighting. As opposed to the pose and expression vectors, which are mainly used to give accurate matches to an input image, our lighting based feature matching is meant to ensure variety and diversity in lighting conditions in order to counteract bias in our datasets. Since the tip of the nose is highly rigid and contains a wide range of surface normal vectors, it well approximates environmental lighting conditions. Thus, motivated by the use of chrome spheres [31], we use a tight crop around the nose as a rough representation of the environmental lighting. The crop is rasterized to a 4×3 grid of RGB values, creating a feature vector $f_{lighting}(I) \in \mathbb{R}^{12 \times 3}$. See Figure 4.



Figure 4. Image I , tight nose crop of I , and rasterized crop that gets flattened to $f_{lighting}(I)$ (three examples are shown). Note that shadows (middle example) and specular highlights (rightmost example) are picked up in the rasterized features.

3.1. k-d Tree Construction

Given an image collection $\Omega = \{I\}$, one can evaluate a feature function $f(I)$ on Ω and use k-means clustering to partition Ω (or a subset of Ω) into k clusters. We use the k-d tree algorithm to recursively build a hierarchy of such clusters, specifying feature functions $f_1(I), \dots, f_n(I)$. For the first $n - 1$ levels, tree nodes at level i subdivide sets of images into a pre-specified number (k_i) of clusters, and the n -th level is flat.

As a specific example, Ω might typically contain a few thousand images. Then, choosing $f_1 = f_{pose}(I)$, $f_2 = f_{lighting}(I)$, $f_3 = f_{expression}(I)$, $k_1 = 9$, and $k_2 = 3$, we would build the cluster hierarchy as follows: first, we extract $f_{pose}(I)$ for all $I \in \Omega$, and split Ω into $k_1 = 9$ clusters based on $f_{pose}(I)$. Within each of the nine pose clusters, we then further split the images into $k_2 = 3$ subclusters based on $f_{lighting}(I)$. The resulting k-d tree is thus defined by (and stored as) nine cluster centers in the pose feature space, 27 cluster centers in the lighting feature space, and a flat list of expression features $f_{expression}(I)$ in each of the 27 lighting clusters.

After building the k-d tree, we can then efficiently retrieve a subset of images by specifying an input image I_{query} , evaluating $f_1(I_{query}), \dots, f_n(I_{query})$, and specifying the number of desired matches at each level of the k-d tree: q_1, \dots, q_n . For example, a query with $q_1 = 2$, $q_2 = 3$, $q_3 = 1$ would select $q_1 = 2$ best matches out of the nine pose feature cluster centers at the first level; then, within each of those 2 pose cluster matches, the query would select all $q_2 = 3$ lighting feature clusters at the second level (to

provide good variation in lighting). Finally, within each of those six image partitions, we'd then pick the closest expression match ($q_3 = 1$), resulting in the retrieval of six images.

Figure 5 shows some results across different datasets, with k-d trees of depth 3, f_1, f_2, f_3 =pose, lighting, expression, $k_1, k_2 = 9, 3$, and $q_1, q_2, q_3 = 2, 3, 1$. Our perceptual features were designed to be as agnostic to facial appearance as possible so that given a hierarchical cluster of images of person A, one can query using images of person B.



Figure 5. Query image on the left, followed by 6 matches from a Tom Cruise dataset, optimized to find nearby pose and expression while maintaining variance in lighting (two examples are shown). The Tom Cruise images were gathered from four separate youtube interviews.

4. Bridging the Domain Gap

Typically, when using synthetic data, one would strive to make the images as realistic-looking as possible in order to minimize the domain gap between the synthetically generated and the in-the-wild images. In contrast, we use a minimally complex rendering pipeline (diffuse shading, simplistic textures, etc.) and bridge the domain gap with unsupervised conditional autoencoders as used in deepfake technology. Using only a small set of curated training data for each subject, we utilize deepfake technology to correlate and transform between synthetic and in-the-wild images. Serendipitously, this allows us to aggressively “annotate” synthetic renders, and subsequently transport that annotation to in-the-wild images (see e.g. Section 4.1).

For the deepfake component of our pipeline, we use a single encoder (E) and fully decoupled dual decoder (D_A, D_B) architecture [49]. During training, aligned crops of both synthetically-generated and in-the-wild images are fed into the joint encoder E , before being passed through the two separate decoders D_A, D_B . An image reconstruction loss

$$f_{recon} = \|D_A(E(I_A)) - I_A\| + \|D_B(E(I_B)) - I_B\|$$

and discriminators G_{image}, G_{latent} paired with asymmetric GAN loss

$$f_{image} = \|G_{image}(D_A(E(I_A)))\| + \|\vec{1} - G_{image}(I_A)\|$$

$$f_{latent} = \|G_{latent}(E(I_B))\| + \|\vec{1} - G_{latent}(E(I_A))\|$$

serve as the main penalty constraints of the decoder outputs during training. As is typical, G_{image} attempts to discern

whether the image comes from Ω_A or the decoder D_A (minimizing f_{image}), whereas D_A attempts to fool the discriminator (maximizing f_{image}). G_{latent} attempts to discern whether the latent image encoding $E(I)$ was generated from an image $I \in \Omega_A$ or $I \in \Omega_B$ (minimizing f_{latent}), whereas E attempts to fool the discriminator (maximizing f_{latent}). The asymmetric GAN loss f_{image} could be easily reflected and made symmetric by incorporating a matching set of discriminators and losses, but in practice, we care specifically about the results of one specific decoder rather than both (e.g. the ‘real’ decoder for appearance capture in Section 5, and the ‘synthetic’ decoder for motion capture in Sections 8 and 9.) The final results of unsupervised training on two sets of facial images $\Omega_A = \{I_A\}$ and $\Omega_B = \{I_B\}$ are a learned encoder E that encodes and correlates Ω_A and Ω_B into a joint embedding space, and two separate decoders D_A, D_B that can decode embedded vectors into images that seemingly belong to either Ω_A or Ω_B (see Figure 6).



Figure 6. Column 1 shows three different examples of an image I_A from a set Ω_A (Tom Cruise data, synthetic renders, synthetic segmentation renders). Column 2 shows $D_A(E(I_A))$. Column 3 and 4 show similar results for I_B and $D_B(E(I_B))$ (Lupita Nyong'o, Tom Cruise, Tom Cruise). Finally, column 5 shows $D_A(E(I_B))$.

4.1. 2D Motion Tracking

The goal of 2D image based tracking is to draw feature points (and/or curves) on images in order to annotate motion. The most successful and robust landmark detection methods typically use CNN based approaches, even though CNNs work with areas of pixels, not codimension one features. The common method of addressing this discrepancy is to leverage area based heatmaps by detecting local extrema to determine tracked points. By drawing on ideas from levelset methods (see e.g. [48]), we instead embed codimension one points and codimension two curves as the boundaries of tracked areas. It is well known that one can track levelsets of functions far more accurately and robustly than extrema.

To enable this approach, we segment the 3-dimensional face surface into regions, with boundaries that are a superset of the codimension one and codimension two features we

wish to track. Then, given a set of in-the-wild images Ω_R , we construct a synthetically rendered dataset Ω_S using the segmented 3D face. Once the deepfake network is properly trained, any image $I_R \in \Omega_R$ can be passed into the encoder and decoded as $D_S(E(I_R))$. Identifying the boundaries between regions on $D_S(E(I_R))$ gives the desired codimension 1 and 2 features that we wish to track. (See Figure 7.)



Figure 7. We train a deepfake network on webcam footage Ω_R and synthetic segmented renders Ω_S . From left to right: real image $I_R \in \Omega_R$, deepfake results $D_S(E(I_R))$, and the boundary contours of $D_S(E(I_R))$ overlaid onto I_R .

5. Appearance Capture

Revolving (theater) stages, colloquially known as turntables, are an effective way of systematically evaluating (and iterating upon) the design of 3D models. Examples can be seen in advertising, car shows, art displays, sculptures for practical special effects, etc. Synthetic renders of steadily rotating computer-generated objects give a holistic view of an object’s geometry and texture, allowing one to avoid 3D geometry that only looks good from a subset of viewpoints. In addition, multi-view stereo synthesis of geometry is well-established in computer vision, and core to many classical methods (visual hull, structure from motion, etc.). Multi-view stereo approaches for facial capture traditionally depend on having fine-grained control over the capture environment, making such approaches impractical for democratization. In contrast, when working with synthetic renders instead of real-world capture data, one has full control over both geometry and camera, which we can leverage to create a dense and precise multi-view pipeline. Our ap-

pearance capture pipeline is thus centered around turntable based evaluation and synthesis of 3D geometry.

5.1. Geometry Estimation

To bootstrap our approach, we utilize an existing single-view based method [53] to get an initial estimate of the geometry. Although such approaches allow for the democratization of appearance capture, they suffer significantly from dataset bias because single-view reconstruction is a highly underconstrained problem. Thus, we prefer a personalized approach when possible (e.g. a classical high-end scanner, a multi-view neural network, depth range, etc.). More specifically, we recommend starting with an estimate of the geometry and subsequently refining it using inverse rendering. We have had particular success using the segmented face texture from Section 4.1 (see Section 8 for more details). Regardless, any method suited for democratization will incur errors in the reconstructed geometry; thus, it is necessary to utilize a texture acquisition method that is able to cross the domain gap between imperfect geometry and real-world images.

5.2. Data Curation

Given estimated geometry, we generate a set of 200 synthetically rendered images with varying pitch and yaw (see the discussion of pose in Section 3). Additionally, given in-the-wild footage of the subject under consideration, we build a k-d tree as discussed in Section 3. Then, for each of the 200 synthetic turntable renders, we use the k-d tree to find a few best matches in pose and expression. The end result is about 200-300 unique images from the in-the-wild footage, all with a relatively fixed expression (matching the synthetic geometry) across a wide range of views. This data curation step enables our method to work on input data ranging from five minute youtube clips of people talking and gesticulating to 30 second shaky self-recorded videos taken on handheld smartphones. See Figure 8.



Figure 8. Example matches to synthetic renders found in real-world video footage using hierarchical clustering. Note that instead of having a fixed pitch and varying yaw as is standard for turntables, we found varying the pitch of the face to be quite useful as well.

5.3. Deepfake Training

Given the 200 synthetic renders and the similarly-sized pruned in-the-wild dataset, we train a deepfake network (see Section 4) with input and output sizes of 384×384 pixels. With these small and heavily correlated datasets, it takes less than an hour for the network model to converge using a

single NVIDIA RTX GPU. By utilizing a slightly-pretrained model as a warm start, one can further cut this down to under half an hour.

5.4. Texture Acquisition: Face

We take a subset of 20 synthetic turntable renders, and for each render I , we generate a pixel-by-pixel rasterization $U(I)$ of the 3D model’s UV coordinates as well as a deepfake result $D(I) = D_R(E(I))$. See Figure 9. Motivated by

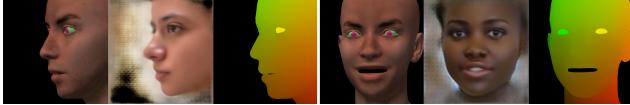


Figure 9. Left to right (two examples shown): Given a synthetic turntable render I , the deepfake result $D(I)$ is well aligned to it, allowing us to generate correspondences between the deepfake pixels and the rasterized UV coordinates $U(I)$ of the model.

photon gathering [27], for each valid pixel, we copy its color $\vec{c} \in D(I)$ into the UV space as a particle sample at its corresponding location $\vec{p} \in U(I)$. After consolidating particle samples from all 20 deepfake images into a single photon map, we use photon gathering to generate the final texture image T . As is typical in photon mapping, the size of the 2D disc around each texel coordinate is increased or decreased to collect a fixed number of particle samples, and the colors of the collected samples are subsequently averaged to obtain the final texel color. We use an inverse-distance weighted average to aggregate the samples. Additionally, each texel is assigned a confidence score inversely proportional to the disc radius. See Figure 10.

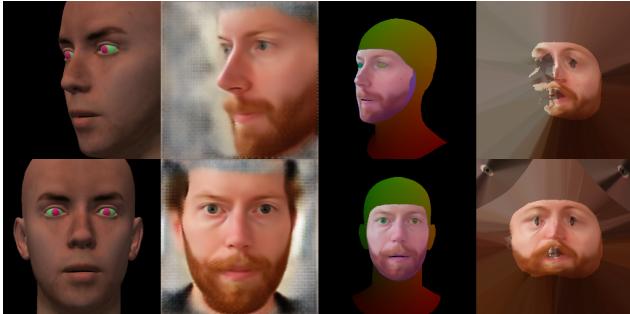


Figure 10. Left to right: Synthetic render, deepfake, high confidence deepfake pixels overlaid onto UV coordinates, and the texture map generated from correspondences between deepfake pixel and UV coordinate pairs. The first row shows results generated using a single view. The second row shows cumulative results generated using multiple views (including the views shown.)

5.5. Texture Acquisition: Head

Deepfake networks are most robust when constrained to a tight crop of the face. In order to ascertain texture in regions where we do not have accurate deepfake pixels, we

use hierarchical clustering (see Section 3) to obtain in-the-wild image matches to the synthetic turntable renders. The matches are close but not perfect, often due to slight variation in expression as well as errors in the geometry estimation step; therefore, we strongly prefer to use aligned deepfake pixels (Section 5.4) whenever possible (see Figure 13). For each synthetic render I , after finding the best in-the-wild image match, we adjust the 3D model to fit that match as closely as possible (as discussed regarding the pose feature in Section 3). Then, we generate $U(I)$ and use the selected in-the-wild image in place of $D(I)$ to generate a texture map along the lines of Section 5.4. See Figure 10. Afterwards, the results are combined with the texture map from Section 5.4, using the per-texel confidence scores mentioned in that section. See Figure 11.

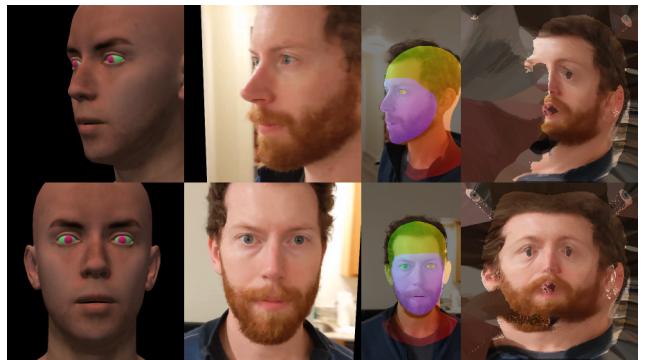


Figure 11. Left to right: Synthetic render, closest dataset match from hierarchical clustering, UV coordinates of the (pose-fitted) synthetic model overlaid onto the source plate of the closest match, and the texture map generated from correspondences between source plate pixel and UV coordinate pairs. The first row shows results generated using a single view. The second row shows cumulative results generated using multiple views (including the views shown).

5.6. Results

We evaluated our pipeline on data ranging from 20 seconds of video from a handheld smartphone camera to 5 minutes of video from youtube celebrity interviews. See Figure 12. Starting from the extraction of frame-by-frame data from raw in-the-wild video input, texture maps for each subject were generated within an hour. Compared to traditional multi-view stereo, our deepfake approach produces fewer misalignment artifacts when paired with imperfect geometry. In particular, the improvements are most striking on non-Caucasian female subjects (see Figure 13), since the monolithic network we use for predicting initial geometry struggles the most on these subjects (presumably due to the fact that they are not well-represented in the dataset used to train the network).

The biggest limitation of our work currently lies with the resolution constraints of deepfake networks. Since our

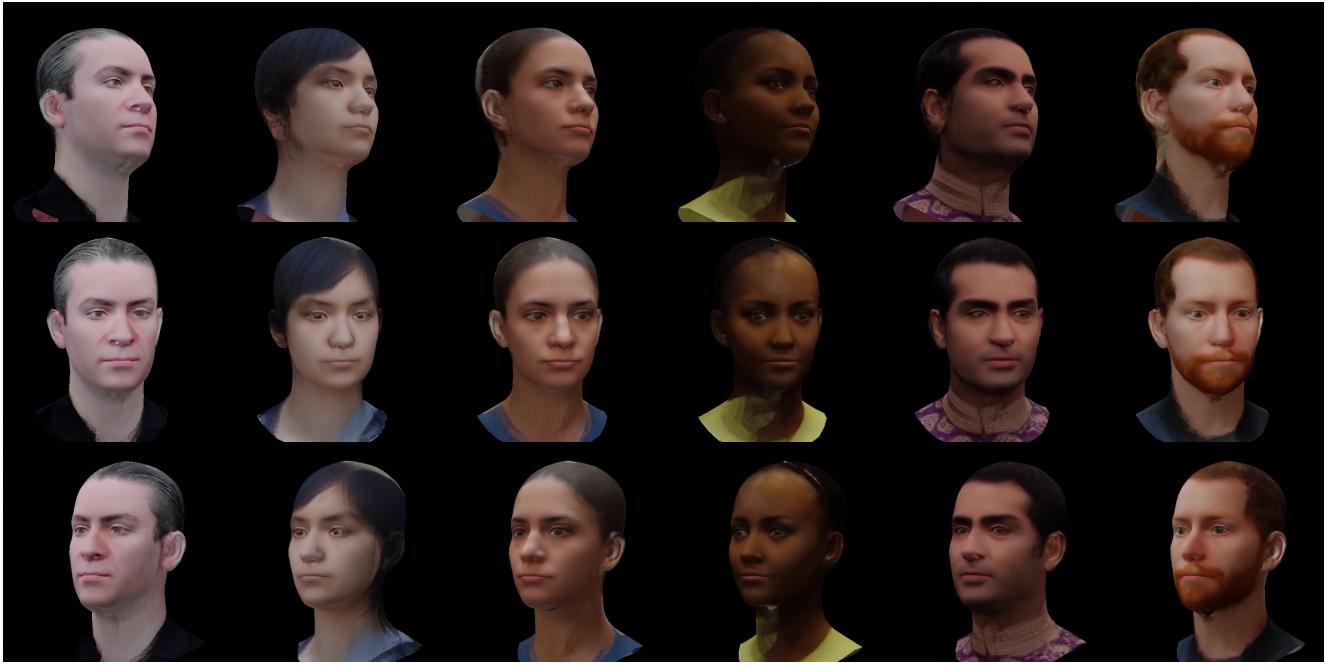


Figure 12. Results on six subjects, rendered from several views. We present the raw results without any manual cleanup in order to illustrate the potential for democratizing our approach. It is straightforward to utilize matting networks, etc. to clean up the texture maps.

approach does not require modifications of existing dual auto-encoder technology, we foresee that the efficacy of our results will improve as deepfake technology becomes scalable to higher resolutions. Additionally, while an automated pipeline that can run in under an hour is already tractable for democratization, we foresee that more aggressive optimizations would be required in order to scale the approach to millions of casual users on their home devices.



Figure 13. The top row shows a texture generated by splatting in-the-wild imagery onto imperfect predicted geometry, with our texture results on the bottom row. Note that the misalignment in the top row is particularly troublesome around the eyes and nose.

An interesting line of future research would be to build texture maps that can be used with more sophisticated light-

ing and shading, such as separated specular, diffuse, and subsurface-scattering texture maps. The turntable synthetic renders would then also need to incorporate lighting into consideration, necessitating a virtual lightstage that achieves similar results to [41] through deepfake technology.

6. Motion Capture

We can transfer any animation rig to the textured geometry resulting from Section 5, e.g. via a volumetric morph defined by surface-curve based boundary constraints [15]. In particular, we begin by transferring the Metahuman [21] joint based rig to our Section 5 results. See Figure 14. The

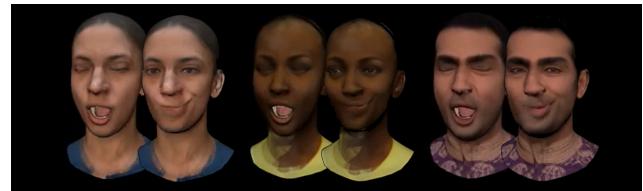


Figure 14. After morphing the metahuman rig to our textured geometry, we can evaluate our appearance capture results on a range of (retargeted) artist directed animations.

Metahuman rig is composed of 167 user controls, 253 shapes, and 683 shape correctives on a surface mesh with roughly 24,000 vertices. For the purposes of motion capture, we simplify this rig by converting it to a linear blendshape rig with jaw joint skinning, discarding the correctives. More in-depth discussions of expression rigs can be found in [20].

For each shape S_i in the original rig, we split it into deformations caused by the jaw rotation and translation (R_i and T_i) and the "de-skinned" offset D_i from the jaw-induced deformation. Given the per-vertex jaw-skinning coefficient tensor M and neutral shape N , each shape can then be described as

$$\begin{aligned} S_i &= M(R_i(N + D_i) + T_i) + (1 - M)(N + D_i) \\ &= (MR_i + (1 - M)I)(N + D_i) + MT_i \end{aligned}$$

where the entries of M vary between 0 and 1 (equal to 0 for vertices only affected by the skull, and 1 for vertices only affected by the jaw). This rig is parametrized by $\vec{w} \in [0, 1]^{253}$ via

$$S(\vec{w}) = (MR(\vec{w}) + (1 - M)I)(N + \sum w_i D_i) + MT(\vec{w})$$

where the jaw rotation matrix $R(\vec{w})$ is generated via Euler angle interpolation of the R_i , and $T(\vec{w}) = \sum w_i T_i$. The nonrigid deformations $S(\vec{w})$ combined with the rigid transformation of the skull P_{rigid} fully defines the range of motion we wish to capture. See Figure 15. Motion tracking is thus defined as extracting the pose P_{rigid} and control rig parameters \vec{w} from images.



Figure 15. From left to right, the neutral face is shown in a rigid pose $P_{rigid}N$ (rotated to the left), with linear deformations $P_{rigid}(N + \sum w_i D_i)$, and with both linear deformations and jaw skinning $P_{rigid}S(\vec{w})$.

In Section 7, we discuss data curation with a specific focus on generating datasets to train deepfakes that can be used to facilitate motion-tracking. In sections 8 and 9, we present two approaches (inverse rendering and learned regression) for deepfake-guided tracking of personalized models.

7. Motion-Tracking Deepfakes

In order to train deepfake models for turntable synthesis, we pruned in-the-wild data to match the synthetically-generated dataset (see Section 3). Conversely, here we augment the synthetic dataset to match the in-the-wild data we wish to track. Ideally, the synthetic images would be generated from a control rig parameter distribution that closely matches the motion present in the in-the-wild data, but acquiring such a distribution is a "chicken-and-egg" problem as one cannot generate such synthetic data until the in-the-wild

footage has already been properly tracked. The problem is doubly difficult when working with a high fidelity rig (as opposed to a simpler PCA rig), since the high-dimensionality of such a rig provides for increased expressivity at the cost of a large nonrealistic-expression subspace.

Contemporaneously with our work, [44] similarly used deepfake technology for motion tracking, addressing the data distribution problem by simultaneously generating animation rigs and control parameter samples via dense 3D motion capture (utilizing head mounted cameras and a pre-defined corpus); however, this approach does not seem to be intended for democratization. With personalized rigs in mind, we propose a dataset generation method that is minimally dependent on the underlying rig and does not require any (high-end) capture setup.

For the sake of exposition, we illustrate our pipeline via a typical example: first, we generate 4,000 random samples by uniformly sampling each control parameter between 0 and 1 subject to some sparsity constraints for each subregion of the face (e.g. no more than five lip shapes would be activated at a time). Each sample is then paired with a rigid pose sampled from a truncated Gaussian distribution (pitch $\in [-10^\circ, 10^\circ]$ and yaw $\in [-80^\circ, 80^\circ]$). Additionally, an in-the-wild dataset is constructed from 3 separate youtube clips of a specific actor. Given both the synthetic and the in-the-wild datasets, a deepfake network is trained to obtain a joint latent-space embedding. Section 7.1 describes how we use latent space analysis to remove part (about 20 percent) of the synthetic dataset, and Section 7.2 describes how we generate solver-bootstrapped synthetic data to augment the synthetic dataset (to obtain approximately 5,000 images in total).

7.1. Dataset Contraction

The increased expressivity of high-dimensional rigs also (unfortunately) leads to unnatural geometric deformations for a wide range of control parameters. Naively sampling the control parameter space results in many expressions that would never be present in in-the-wild data. Thus, we prune randomly generated samples via latent space analysis of trained deepfakes. To accomplish this, we first train an unsupervised deepfake model using the randomly-generated synthetic and in-the-wild datasets. Then, taking advantage of the unsupervised clustering that occurs as a byproduct of training the deepfake network, we identify (and delete) *out-of-distribution* synthetic data by finding the outliers (as compared to in-the-wild data) in the deepfake network's latent space. See Figure 16.

7.2. Dataset Expansion

Similar to our dataset contraction approach, we perform dataset expansion by (instead) finding outliers in the real data (as compared to synthetic data), characterizing them as *gaps* in the synthetic dataset, and subsequently utilizing existing



Figure 16. Synthetic renders I_S (left column), and the closest in-the-wild images I_R (middle column) as measured in the latent space. The top row shows an example of ‘out of distribution’ synthetic data, contrasted with the bottom row which shows ‘in distribution’ synthetic data. Comparing $D_S(E(I_R))$ (right column) to I_S (left column) gives one a sense of the distance between $E(I_S)$ and $E(I_R)$.

motion capture solvers to bootstrap *gap-filling* synthetic data. See Figure 17. Extracting control parameters from real images is precisely the problem we wish to solve, so we are again faced with a chicken-and-egg problem. Solely relying on existing solvers to expand the synthetic dataset is not only over-reliant on the accuracy of the solvers, but also causes the network to overfit to and mimic these solvers (as opposed to being trained to match the in-the-wild footage in an optimal way). To avoid this, we only use off-the-shelf solvers to bootstrap the dataset expansion, and subsequently augment bootstrapped data with random jittering and interpolation. More specifically, given a control parameter sample \vec{w} generated from an existing solver, we add per-control random noise to generate new samples $D\vec{w}$, D being a diagonal matrix with diagonal entries $d_i \in [0.8, 1.2]$. In addition, given two control parameter samples \vec{w}_1 and \vec{w}_2 , we interpolate between them via $\alpha\vec{w}_1 + (1 - \alpha)\vec{w}_2$ with $\alpha \in (0, 1)$. Of course, any newly added samples can be examined for outliers (and pruned if necessary) using the approach in Section 7.1.

8. Inverse rendering

Inverse rendering differentiates through the rendering pipeline in order to optimize parameters for geometry, material properties, lighting, etc. The restriction that the renderer has to be fully differentiable limits its visual fidelity, and thus its ability to match photorealistic in-the-wild images. Most practitioners are well-aware of such limitations. Here, we present a promising approach for using inverse rendering on in-the-wild images, by utilizing segmentation deepfakes (as described in Section 4.1). Inferencing a segmented texture



Figure 17. In-the-wild images I_R (left column) and the closest synthetic renders I_S (middle column) as measured in the latent space. The top row shows ‘out of distribution’ in-the-wild data, contrasted with the bottom row which shows ‘in distribution’ in-the-wild data. Comparing $D_S(E(I_R))$ (right column) to I_S (middle column) gives one a sense of the distance between $E(I_R)$ and $E(I_S)$.

from an in-the-wild image (using a trained deepfake) allows one to more robustly determine parameters via an inverse rendering approach.

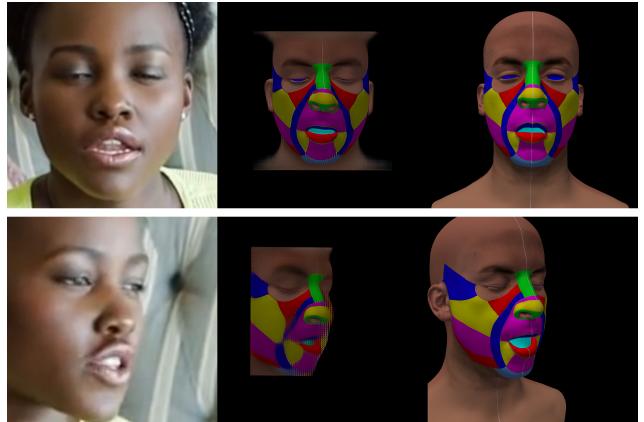


Figure 18. Two examples of inverse rendering solves. From left to right: input frame I_R , $D_S(E(I_R))$ rigidly transformed to the blendshape rig’s image coordinate frame, and a rendering of the blendshape rig with solved parameters.

To do this, we first implement our blendshape rig as a fully-differentiable pytorch3d [51] module. Given a fixed segmentation texture T_{seg} (as opposed to a photorealistic face texture T_{real}),

$$F_{seg}(p, y, \vec{w}) = F(P_{rigid}(p, y)S(\vec{w}); T_{seg})$$

is a differentiable image generating function parametrized by (p, y, \vec{w}) . Then, given an in-the-wild dataset, we use the method described in Section 7 to generate a synthetic

segmentation dataset (Ω_S) that matches the statistical distribution of the in-the-wild data (Ω_R); afterwards, a deepfake model is trained between the two datasets.

Given an in-the-wild image I_R , we solve

$$\arg \min_{p,y,\vec{w}} \|F_{seg}(p,y,\vec{w}) - D_S(E(I_R))\|.$$

using a pixel-wise L2 distance and block coordinate descent, i.e. each epoch first considers the rigid pose, then the subset of control parameters active for jaw-skinning, finally followed by the rest of the control parameters. See Figures 18 and 19.



Figure 19. A breakdown of the block gradient descent algorithm. Given the deepfake result on the left, we iteratively solve for pose, jaw, and expression parameters (from left to right).

For the sake of comparison, we experimented with both realistic and segmented textures, training separate deepfake models to inference either the segmented or realistic texture onto in-the-wild images. We observed lower sensitivity to the initial state and faster convergence when using the segmented (as opposed to the realistic) texture. See Figure 20. In particular, one can typically use the front-facing neutral expression as a robust starting point for every frame; however, starting from an initially computed pose (as described in Section 3) further increases robustness. Note that subspace analysis (described in Section 7.1) can also be used to provide a robust initial guess.

In our experiments, we were able to solve for full in-the-wild sequences without manual supervision at 10-30 seconds per frame. Moreover, the ability to use the neutral expression as an initial guess for every frame allows for parallelization. This scalable offline-solver can be used to generate ground truth datapoints for training real-time regressor models (see Section 9).

8.1. Lip-Sync Deepfakes

One well-known limitation of most motion tracking pipelines is that they struggle to capture the subtle lip/mouth movements required for a high fidelity sync to an audio track. To address this issue, we utilize a specialized neural network trained to inference images of the lip/mouth from speech signal input. Specifically, we use the pretrained Wav2Lip model [50], which can be used to obtain 2D image animations of lip/mouth movements given only an audio clip and a single reference image of a (realistically-textured) synthetic face in the neutral pose. One advantage of using such 2D lipsync models is that the data required to train them, i.e.

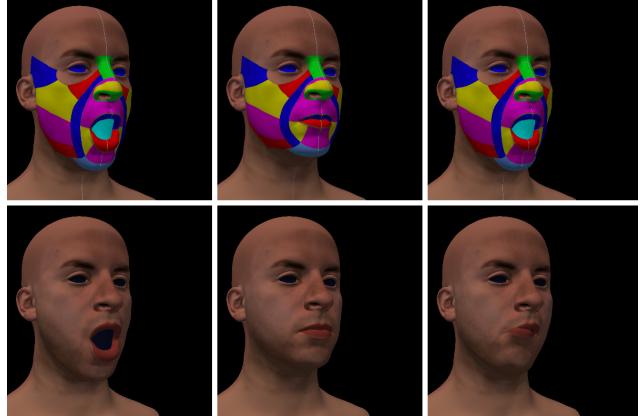


Figure 20. A simple test comparing our use of $F_{seg}(p,y,\vec{w})$ (top row) with the use of a more photorealistic $F_{real}(p,y,\vec{w}) = F(P_{rigid}(p,y)S(\vec{w}); T_{real})$ (bottom row) for the inverse rendering. Given the same target pose and expression p^*, y^*, \vec{w}^* (first column), the inverse renderer using T_{seg} starting from the correct rigid pose (column 2) converges as expected (column 3); replacing T_{seg} with T_{real} gives poor results.

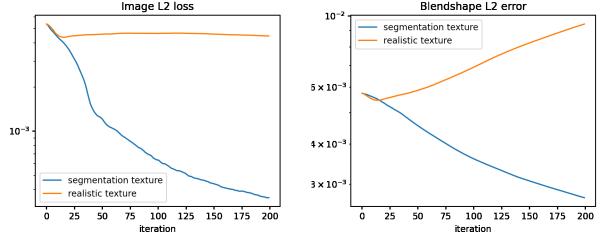


Figure 21. Plots showing the L2 image loss (which our inverse renderer optimizes over), as well as the L2 error on \vec{w} . The blue line corresponds to the first row of Figure 20 using F_{seg} , and the orange line corresponds to row 2 using F_{real} .

videos of people talking, is much more commonly available than data for 3D talking heads. Once again, it is worth noting that we leverage an off-the-shelf pretrained neural network, and thus are able to obtain improved results as research on such networks progresses.

Following this approach leads to two separate videos: our segmented deepfake obtained from inferencing the source video input, and a frontal video of lip/mouth motion obtained from inferencing the source audio input. Notably, we use inverse rendering separately on each video, and subsequently blend the two sets of animation parameters together. This allows one to focus on the rigid pose and expression parameters in one video without worrying too much about the fidelity of the subtle lip/mouth motions, while focusing on highly detailed lip/mouth motion in the other video with fixed pose and (non-lip/mouth) expression (i.e. expressionless, and in the front-facing neutral pose). Essentially, we decompose a difficult problem into two less difficult components, which can each benefit from known constraints and

importance metrics. See Figure 22.

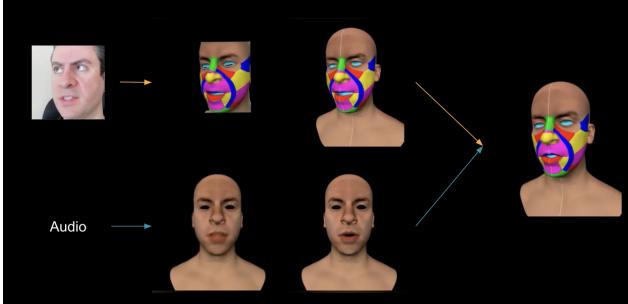


Figure 22. Our proposed approach to enhance subtle lip/mouth movements. In the first row, we generate our segmented deepfake video (obtained from inferencing the source video input) and run inverse rendering to obtain control parameters. In the second row, we generate a frontal video of lip/mouth motion (obtained from inferencing the source audio input) and again run inverse rendering to obtain control parameters. Finally, we blend parameters for rigid pose and (non-lip/mouth) expressions from the first row with parameters for lip/mouth motion from the second row.

9. Learned regression

In this Section, we present work on using the learned encoder E of a deepfake network in order to create a personalized regression model for motion tracking. As discussed in Section 7, we first create a synthetic dataset that well-mimics the in-the-wild data we wish to track (as measured in the deepfake latent space). Then, given ground truth synthetic parameters $p^{(i)}, y^{(i)}, \vec{w}^{(i)}$ paired with generated images $I_S = F_{seg}(p^{(i)}, y^{(i)}, \vec{w}^{(i)})$, a straightforward approach would be to train a regressor that goes from the latent embeddings of synthetic renders $E(I_S)$ to ground truth $p^{(i)}, y^{(i)}, \vec{w}^{(i)}$; subsequently, that regressor could be inferenced on real-image embeddings $E(I_R)$ in order to estimate motion parameters for I_R .

In order to reduce the domain gap even further, we project both the synthetic data and the in-the-wild images into a ‘synthetic’ embedding space, training and inferencing on $\vec{e} = E(D_S(E(I_S)))$ and $\vec{e} = E(D_S(E(I_R)))$ instead of $E(I_S)$ and $E(I_R)$. Moreover, in order to incorporate in-the-wild images during training, we use a pretrained landmark detection network \hat{L} (see [10]) to generate sparse landmarks $\vec{m} = \hat{L}(D_S(E(I_S)))$ and $\vec{m} = \hat{L}(D_S(E(I_R)))$ for weak supervision (used only during training). With pose denoted as $\vec{p} = (p, y)$, the synthetic dataset Ω_S contains datapoints of the form $(\vec{e}, \vec{p}, \vec{w}, \vec{m})$, while the in-the-wild dataset Ω_R contains datapoints of the form (\vec{e}, \vec{m}) lacking ground truth pose and control parameters.

We designed our network model to mimic traditional solvers, with modules *sequentially* predicting pose \vec{p} , jaw control parameters \vec{w}_j (a subset of \vec{w}), and the remaining control parameters \vec{w}_c (which we will refer to as the expression controls).

Given a latent embedding \vec{e} , we denote the pose network as $P(\vec{e})$, the jaw network as $J(\vec{e}, P(\vec{e}))$, the expression network as $W(\vec{e}, P(\vec{e}), J(\vec{e}))$, and our landmark prediction network as $L(P(\vec{e}), J(\vec{e}), W(\vec{e}))$. P, J, W, L are all multi-layer perceptron networks, composed of multiple “fully-connected / leaky ReLU” stacks with a sigmoid activation layer. Our network model is trained in three stages as illustrated in Figure 23.

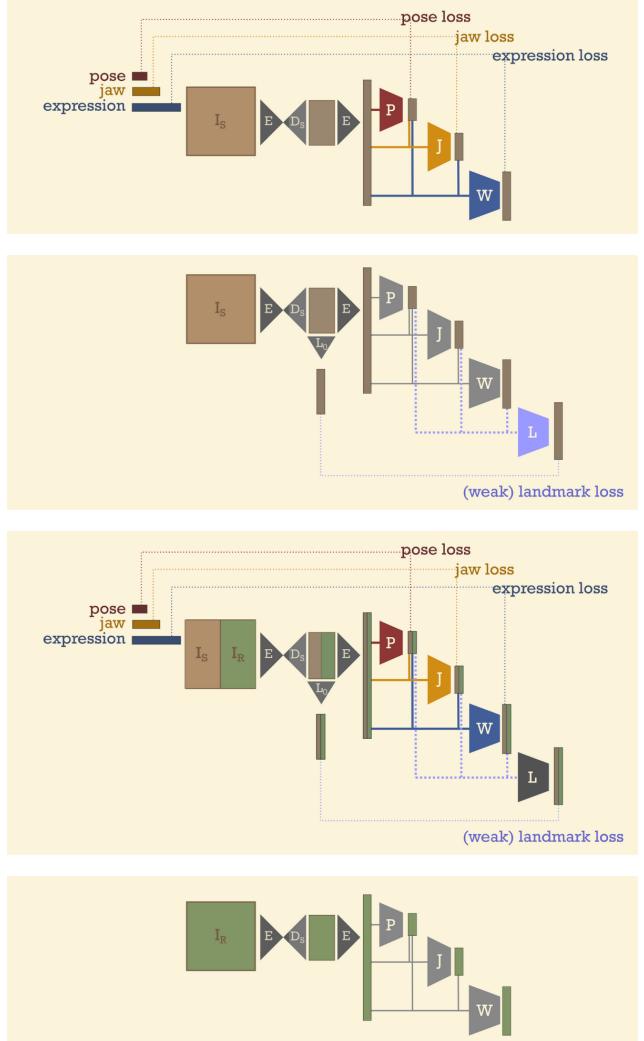


Figure 23. Diagram of our regressor architecture, with the three steps of training shown top to bottom (followed by the evaluation step). Grey denotes frozen layers. The pose, jaw, and expression losses are driven only by the synthetic data (shown in brown). The in-the-wild data (shown in green) is incorporated purely through weak supervision, achieved via landmarks.

9.1. Training

In the first stage, P, J , and W are trained using only synthetically generated data with known ground truth (Figure 23, top). $P(\vec{e})$ learns to predict \vec{p} directly from \vec{e} , $J(\vec{e}, P(\vec{e}))$

learns to predict \vec{w}_j from \vec{e} and $P(\vec{e})$, and $W(\vec{e}, P(\vec{e}), J(\vec{e}))$ learns to predict \vec{w}_c from \vec{e} , $P(\vec{e})$, and $J(\vec{e})$. During training, the loss for J not only compares $J(\vec{e}, P(\vec{e}))$ to \vec{w}_j as expected but also contains an equivalent term using the ground truth \vec{p} in place of $P(\vec{e})$, i.e.

$$\mathcal{L}_J = \sum_{\Omega_S} ||J(\vec{e}, P(\vec{e})) - \vec{w}_j|| + ||J(\vec{e}, \vec{p}) - \vec{w}_j||.$$

Similarly,

$$\mathcal{L}_W = \sum_{\Omega_S} ||W(\vec{e}, P(\vec{e}), J(\vec{e})) - \vec{w}_c|| + ||W(\vec{e}, \vec{p}, \vec{w}_j) - \vec{w}_c||.$$

In the second stage, P, J, W are frozen while we train a minimal landmark network L that learns to loosely predict the estimates $\vec{m} = \hat{L}(D_S(E(I_S)))$ again using only synthetic data (Figure 23, second row). To account for noisy landmarks, which is a known issue in “ground truth” annotations [19] and compounded by our choice of solver (see the discussion in Section 7.2), we use a loss that vanishes when the distance between each marker prediction L_k and its corresponding estimate \vec{m}_k is less than some threshold $\delta > 0$, i.e.

$$\begin{aligned} \mathcal{L}_L = \sum_{\Omega_S} \sum_k \max(0, ||L_k(P(\vec{e}), J(\vec{e}), W(\vec{e})) - \vec{m}_k|| - \delta) \\ + \max(0, ||L_k(\vec{p}, \vec{w}_j, \vec{w}_c) - \vec{m}_k|| - \delta). \end{aligned}$$

In our experiments, we used $\delta = 0.01$ (with ground truth landmarks normalized between 0 and 1).

In the third stage, P, J, W are fine-tuned with L frozen, utilizing both the synthetic and in-the-wild data (Figure 23, third row). Although $\mathcal{L}_P, \mathcal{L}_J, \mathcal{L}_W$ retain their original form, we no longer use the second term in \mathcal{L}_L since it is unavailable on in-the-wild data. One could still include that second term for Ω_S , but we prefer to treat Ω_R and Ω_S more similarly, and modify the landmark loss to be

$$\mathcal{L}_L^* = \sum_{\Omega_S \cup \Omega_R} \sum_k \max(0, ||L_k(P(\vec{e}), J(\vec{e}), W(\vec{e})) - \vec{m}_k|| - \delta).$$

As shown in Figure 24, this third training step is crucial and greatly improves the quality of the results.

9.2. Results and Discussion

Figure 25 shows results on three separate subjects, with no ground truth parameters (obtained or computed). For each example, a deepfake model was trained from scratch using a few thousand in-the-wild images as well as a similarly sized synthetic dataset (see Section 7) at 256×256 image resolution. The resulting deepfake model was then used to extract 512-dimensional latent embeddings from both datasets, and a regressor was trained to predict motion capture parameters from these latent embeddings. Importantly,

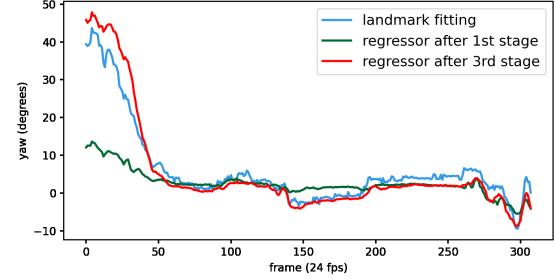


Figure 24. Yaw estimation plotted across a single video of 300+ frames (six sample frames are shown for both the video and the corresponding output from the regressor). The predictions obtained using sparse landmarks (blue) are extremely noisy (as compared to the regressors). The 1st stage regressor (green) trained on only synthetic data struggles to accurately predict yaw on poses that are too far from the neutral, which is not the case for the final regressor (red) trained with the aid of weak supervision via landmarks.



Figure 25. From left to right (three examples shown): input frame I_R , synthetic deepfake $D_S(E(I_R))$, rendering of the 3D model with pose, jaw, and expression regressed from $E(D_S(E(I_R)))$. Interestingly, despite the degradation in image quality of $D_S(E(I_R))$, our regressor still outputs reasonable results.

the low-dimensionality of both the embeddings and the motion capture parameters enables quite efficient training of the regressor (only requiring minutes on a single GPU). The bot-

tleneck is the 6-8 hours required to train a deepfake model, which is something we expect that the deepfake community will improve upon.

The resolution limitations of the deepfake technology similarly limit the fidelity of our results, in particular in the lip region; thus, region-specific encodings (such as the ones described in Section 8.1) are a promising avenue of future research. A deeper investigation into and a comparison of various animation rig parametrizations would also be interesting. Finally, although we propose an alternative to large models and large scale data collection, a hybrid approach that leverages pretrained models might help to alleviate the training time needed for the deepfake networks.

10. Summary

In this paper, we presented various techniques that utilize unsupervised autoencoder deep-fake neural networks in order to create personalized facial appearance and motion capture pipelines. Notably, only a minimal amount of subject-specific in-the-wild imagery is required, since we were able to leverage synthetically created “ground truth” data during training. Such personalized facial appearance and motion capture pipelines bypass issues with bias that plague large-scale data collection and large pre-trained monolithic models. Our new approach has obvious potential benefits beyond appearance/motion capture and retargeting. For example, the reliance on only subject-specific in-the-wild data could be leveraged to parametrize and analyze the facial motion present in videos in order to detect frame based video forgery (e.g. as created by deepfakes).

11. Acknowledgments

This work was supported by a grant from Meta (née Facebook) and in part by ONR N00014-13-1-0346, ONR N00014-17-1-2174. We would also like to thank Industrial Light and Magic for supporting our initial experiments on deepfake based facial capture, Robert Huang and William Tsu for their kind donation of an NVIDIA Titan X GPU which was used to run experiments, and Epic Games for their help with the Metahuman rig.

References

- [1] M. Bao, M. Cong, S. Grabli, and R. Fedkiw. High-quality face capture using anatomical muscles. *CoRR*, abs/1812.02836, 2018. [2](#)
- [2] T. Beeler, B. Bickel, P. Beardsley, B. Sumner, and M. Gross. High-quality single-shot capture of facial geometry. In *ACM SIGGRAPH 2010 papers*, pages 1–9. 2010. [2](#)
- [3] T. Beeler, M. Gross, P. Gotardo, J. Riviere, and D. Bradley. Medusa facial capture system, 2022. [2](#)
- [4] T. Beeler, F. Hahn, D. Bradley, B. Bickel, P. Beardsley, C. Gotsman, R. W. Sumner, and M. Gross. High-quality passive facial performance capture using anchor frames. In *ACM SIGGRAPH 2011 papers*, pages 1–10. 2011. [2](#)
- [5] K. S. Bhat, R. Goldenthal, Y. Ye, R. Mallet, and M. Kopferwas. High fidelity facial animation capture and retargeting with contours. In *Proceedings of the 12th ACM SIGGRAPH/eurographics symposium on computer animation*, pages 7–14, 2013. [2](#)
- [6] K. S. Bhat, R. Goldenthal, Y. Ye, R. Mallet, and M. Kopferwas. High fidelity facial animation capture and retargeting with contours. In *Proceedings of the 12th ACM SIGGRAPH/Eurographics Symposium on Computer Animation, SCA ’13*, page 7–14, New York, NY, USA, 2013. Association for Computing Machinery. [2](#)
- [7] V. Blanz and T. Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, pages 187–194, 1999. [2](#)
- [8] J. Booth, A. Roussos, S. Zafeiriou, A. Ponniah, and D. Dunaway. A 3d morphable model learnt from 10,000 faces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5543–5552, 2016. [2](#)
- [9] C. Bregler. Motion capture technology for entertainment [in the spotlight]. *IEEE Signal Processing Magazine*, 24(6):160–158, 2007. [2](#)
- [10] A. Bulat and G. Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem?(and a dataset of 230,000 3d facial landmarks). In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1021–1030, 2017. [3, 12](#)
- [11] J. Buolamwini and T. Gebru. Gender shades: Inter-sectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91. PMLR, 2018. [2](#)
- [12] D. Cheng, Y. Gong, S. Zhou, J. Wang, and N. Zheng. Person re-identification by multi-channel parts-based cnn with improved triplet loss function. In *Proceedings of the iEEE conference on computer vision and pattern recognition*, pages 1335–1344, 2016. [4](#)
- [13] S. Cheng, M. M. Bronstein, Y. Zhou, I. Kotsia, M. Pantic, and S. Zafeiriou. Meshgan: Non-linear 3d morphable models of faces. *CoRR*, abs/1903.10384, 2019. [2](#)
- [14] G. G. Chrysos, E. Antonakos, P. Snape, A. Asthana, and S. Zafeiriou. A comprehensive performance evaluation of deformable face tracking “in-the-wild”. *International Journal of Computer Vision*, 126(2):198–232, 2018. [3](#)
- [15] M. Cong, M. Bao, J. L. E, K. S. Bhat, and R. Fedkiw. Fully automatic generation of anatomical face simulation models. In *Proceedings of the 14th ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, pages 175–183, 2015. [8](#)
- [16] M. Cong, L. Lan, and R. Fedkiw. Local geometric indexing of high resolution data for facial reconstruction from sparse markers. *CoRR*, abs/1903.00119, 2019. [2](#)
- [17] P. Debevec. The light stages and their applications to photo-real digital actors. *SIGGRAPH Asia*, 2(4):1–6, 2012. [2](#)
- [18] P. Debevec, T. Hawkins, C. Tchou, H.-P. Duiker, and W. Sarokin. Acquiring the Reflectance Field of a Human Face. In *SIGGRAPH*, New Orleans, LA, July 2000. [2](#)

- [19] X. Dong, S.-I. Yu, X. Weng, S.-E. Wei, Y. Yang, and Y. Sheikh. Supervision-by-registration: An unsupervised approach to improve the precision of facial landmark detectors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 360–368, 2018. 13
- [20] B. Egger, W. A. Smith, A. Tewari, S. Wuhrer, M. Zollhoefer, T. Beeler, F. Bernard, T. Bolkart, A. Kortylewski, S. Romdhani, et al. 3d morphable face models—past, present, and future. *ACM Transactions on Graphics (TOG)*, 39(5):1–38, 2020. 3, 8
- [21] E. Games. Metahuman creator, 2021. 8
- [22] B. Gecer, S. Ploumpis, I. Kotsia, and S. Zafeiriou. Ganfit: Generative adversarial network fitting for high fidelity 3d face reconstruction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1155–1164, 2019. 3
- [23] A. Ghosh, G. Fyffe, B. Tunwattanapong, J. Busch, X. Yu, and P. Debevec. Multiview face capture using polarized spherical gradient illumination. In *Proceedings of the 2011 SIGGRAPH Asia Conference*, pages 1–10, 2011. 2
- [24] S. Grabli, M. Bao, P. Karelfelt, A. Ferrall-Nunge, J. Yost, R. Fedkiw, C. Phillips, P. Helman, and L. Estebecorena. On-set facial performance capture and transfer to a three-dimensional computer-generated model, U.S. Patent 16681300, July. 2021. 2
- [25] D. Hendler, L. Moser, R. Battulwar, D. Corral, P. Cramer, R. Miller, R. Cloudsdale, and D. Roble. Avengers: capturing thanos’s complex face. In *ACM SIGGRAPH 2018 Talks*, pages 1–2. 2018. 2
- [26] A. E. Ichim, S. Bouaziz, and M. Pauly. Dynamic 3d avatar creation from hand-held video input. *ACM Transactions on Graphics (ToG)*, 34(4):1–14, 2015. 3
- [27] H. W. Jensen. *Realistic image synthesis using photon mapping*. AK Peters/crc Press, 2001. 7
- [28] A. Jourabloo, F. D. la Torre, J. M. Saragih, S. Wei, T. Wang, S. Lombardi, D. Belko, A. Trimble, and H. Badino. Robust egocentric photo-realistic facial expression transfer for virtual reality. *CoRR*, abs/2104.04794, 2021. 3
- [29] Z. Khan and Y. Fu. One label, one billion faces: Usage and consistency of racial categories in computer vision. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pages 587–597, 2021. 2
- [30] J. T. Klosowski, M. Held, J. S. Mitchell, H. Sowizral, and K. Zikan. Efficient collision detection using bounding volume hierarchies of k-dops. *IEEE transactions on Visualization and Computer Graphics*, 4(1):21–36, 1998. 3
- [31] H. Landis. Production-ready global illumination. *Siggraph course notes*, 16(2002):11, 2002. 4
- [32] J. P. Lewis, K. Anjyo, T. Rhee, M. Zhang, F. Pighin, and Z. Deng. Practice and Theory of Blendshape Facial Models. In S. Lefebvre and M. Spagnuolo, editors, *Eurographics 2014 - State of the Art Reports*. The Eurographics Association, 2014. 3
- [33] H. Li, J. Yu, Y. Ye, and C. Bregler. Realtime facial animation with on-the-fly correctives. *ACM Trans. Graph.*, 32(4):42–1, 2013. 3
- [34] T. Li, T. Bolkart, M. J. Black, H. Li, and J. Romero. Learning a model of facial shape and expression from 4d scans. *ACM Trans. Graph.*, 36(6):194–1, 2017. 2
- [35] Z. Li, J.-i. Imai, and M. Kaneko. Facial-component-based bag of words and phog descriptor for facial expression recognition. In *2009 IEEE International Conference on Systems, Man and Cybernetics*, pages 1353–1358. IEEE, 2009. 4
- [36] Y. Liu, A. Jourabloo, W. Ren, and X. Liu. Dense face alignment. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 1619–1628, 2017. 3
- [37] S. Lombardi, J. M. Saragih, T. Simon, and Y. Sheikh. Deep appearance models for face rendering. *CoRR*, abs/1808.00362, 2018. 3
- [38] M. M. Loper and M. J. Black. OpenDR: An approximate differentiable renderer. In *Computer Vision – ECCV 2014*, volume 8695 of *Lecture Notes in Computer Science*, pages 154–169. Springer International Publishing, Sept. 2014. 2
- [39] H. Luo, K. Nagano, H.-W. Kung, Q. Xu, Z. Wang, L. Wei, L. Hu, and H. Li. Normalized avatar synthesis using stylegan and perceptual refinement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11662–11672, June 2021. 3
- [40] S. Ma, T. Simon, J. M. Saragih, D. Wang, Y. Li, F. D. la Torre, and Y. Sheikh. Pixel codec avatars. *CoRR*, abs/2104.04638, 2021. 3
- [41] W.-C. Ma, T. Hawkins, P. Peers, C.-F. Chabert, M. Weiss, and P. Debevec. Rapid acquisition of specular and diffuse normal maps from polarized spherical gradient illumination. In *Proceedings of the 18th Eurographics Conference on Rendering Techniques, EGSR’07*, page 183–194, Goslar, DEU, 2007. Eurographics Association. 8
- [42] W.-C. Ma, A. Jones, T. Hawkins, J.-Y. Chiang, and P. Debevec. A high-resolution geometry capture system for facial performance. In *ACM SIGGRAPH 2008 Talks, SIGGRAPH ’08*, New York, NY, USA, 2008. Association for Computing Machinery. 2
- [43] M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, I. D. Raji, and T. Gebru. Model cards for model reporting. *Proceedings of the Conference on Fairness, Accountability, and Transparency*, Jan 2019. 2
- [44] L. Moser, C. Chien, M. Williams, J. Serra, D. Hendler, and D. Roble. Semi-supervised video-driven facial animation transfer for production. *ACM Transactions on Graphics (TOG)*, 40(6):1–18, 2021. 3, 9
- [45] L. Moser, D. Hendler, and D. Roble. Masquerade: fine-scale details for head-mounted camera motion capture data. In *ACM SIGGRAPH 2017 Talks*, pages 1–2. 2017. 2
- [46] K. Nagano, H. Luo, Z. Wang, J. Seo, J. Xing, L. Hu, L. Wei, and H. Li. Deep face normalization. *ACM Trans. Graph.*, 38(6), nov 2019. 3
- [47] I. Navarro, D. Kneubuehler, T. Verhulsdonck, E. D. Du Bois, W. Welch, V. Verma, I. Sachs, and K. Bhat. Fast facial animation from video. In *ACM SIGGRAPH 2021 Talks*, pages 1–2. 2021. 3
- [48] S. Osher and R. P. Fedkiw. Level set methods: an overview and some recent results. *Journal of Computational physics*, 169(2):463–502, 2001. 5

- [49] I. Perov, D. Gao, N. Chervoni, K. Liu, S. Marangonda, C. Umé, M. Dpfks, C. S. Facenheim, L. RP, J. Jiang, S. Zhang, P. Wu, B. Zhou, and W. Zhang. Deepfacelab: Integrated, flexible and extensible face-swapping framework, 2021. 5
- [50] K. R. Prajwal, R. Mukhopadhyay, V. P. Namboodiri, and C. Jawahar. A lip sync expert is all you need for speech to lip generation in the wild. In *Proceedings of the 28th ACM International Conference on Multimedia*, MM '20, page 484–492, New York, NY, USA, 2020. Association for Computing Machinery. 11
- [51] N. Ravi, J. Reizenstein, D. Novotny, T. Gordon, W.-Y. Lo, J. Johnson, and G. Gkioxari. Accelerating 3d deep learning with pytorch3d. arXiv:2007.08501, 2020. 2, 10
- [52] J. Riviere, P. Gotardo, D. Bradley, A. Ghosh, and T. Beeler. Single-shot high-quality facial geometry and skin appearance capture. *ACM Trans. Graph.*, 39(4), jul 2020. 2
- [53] S. Sanyal, T. Bolkart, H. Feng, and M. J. Black. Learning to regress 3d face shape and expression from an image without 3d supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7763–7772, 2019. 3, 6
- [54] O. Sorkine. Least-squares rigid motion using svd. 4
- [55] R. W. Sumner and J. Popović. Deformation transfer for triangle meshes. *ACM Transactions on graphics (TOG)*, 23(3):399–405, 2004. 2
- [56] Q. Tan, L. Gao, Y. Lai, and S. Xia. Variational autoencoders for deforming 3d mesh models. *CoRR*, abs/1709.04307, 2017. 2, 3
- [57] A. Tewari, M. Zollhoefer, F. Bernard, P. Garrido, H. Kim, P. Perez, and C. Theobalt. High-fidelity monocular face reconstruction based on an unsupervised model-based face autoencoder. *IEEE transactions on pattern analysis and machine intelligence*, 42(2):357–370, 2018. 3
- [58] A. Tewari, M. Zollhöfer, P. Garrido, F. Bernard, H. Kim, P. Pérez, and C. Theobalt. Self-supervised multi-level face model learning for monocular reconstruction at over 250 hz. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 3
- [59] J. Thies, M. Zollhöfer, M. Stamminger, C. Theobalt, and M. Nießner. Facevr: Real-time facial reenactment and eye gaze control in virtual reality. *CoRR*, abs/1610.03151, 2016. 3
- [60] S. Umeyama. Least-squares estimation of transformation parameters between two point patterns. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 13(04):376–380, 1991. 4
- [61] T. Weise, S. Bouaziz, H. Li, and M. Pauly. Realtime performance-based facial animation. *ACM transactions on graphics (TOG)*, 30(4):1–10, 2011. 3
- [62] H. Yang, H. Zhu, Y. Wang, M. Huang, Q. Shen, R. Yang, and X. Cao. Facescape: a large-scale high quality 3d face dataset and detailed riggable 3d face prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 601–610, 2020. 2
- [63] L. Zhang, N. Snavely, B. Curless, and S. M. Seitz. Spacetime faces: High-resolution capture for modeling and animation. In *ACM Annual Conference on Computer Graphics*, pages 548–558, August 2004. 2
- [64] M. Zollhöfer, J. Thies, P. Garrido, D. Bradley, T. Beeler, P. Pérez, M. Stamminger, M. Nießner, and C. Theobalt. State of the art on monocular 3d face reconstruction, tracking, and applications. In *Computer Graphics Forum*, volume 37, pages 523–550. Wiley Online Library, 2018. 2