Democratizing the Creation of Animatable Facial Avatars

YILIN ZHU, Stanford University, USA

DALTON OMENS, Stanford University, USA and Epic Games, USA

HAODI HE, Stanford University, USA

RON FEDKIW, Stanford University, USA and Epic Games, USA

In high-end visual effects pipelines, a customized (and expensive) light stage system is (typically) used to scan an actor in order to acquire both geometry and texture for various expressions. Aiming towards democratization, we propose a novel pipeline for obtaining geometry and texture as well as enough expression information to build a customized person-specific animation rig without using a light stage or any other high-end hardware (or manual cleanup). A key novel idea consists of warping real-world images to align with the geometry of a template avatar and subsequently projecting the warped image into the template avatar's texture; importantly, this allows us to leverage baked-in real-world lighting/texture information in order to create surrogate facial features (and bridge the domain gap) for the sake of geometry reconstruction. Not only can our method be used to obtain a neutral expression geometry and de-lit texture, but it can also be used to improve avatars after they have been imported into an animation system (noting that such imports tend to be lossy, while also hallucinating various features). Since a default animation rig will contain template expressions that do not correctly correspond to those of a particular individual, we use a Simon Says approach to capture various expressions and build a person-specific animation rig (that moves like they do). Our aforementioned warping/projection method has high enough efficacy to reconstruct geometry corresponding to each expressions.

1 INTRODUCTION

The use of personalized avatars has become increasingly prevalent in a wide range of applications: massively multiplayer online games (e.g. Roblox), the entertainment and animation industry (e.g. MetaHumans from Epic Games), virtual/augmented reality experiences (e.g. Apple, Meta, Snapchat), and video conferencing platforms (e.g. Microsoft Teams, Zoom). Personalized avatars have been shown [Sailer et al. 2017] to enhance user engagement and satisfaction in a variety of online settings including gaming, virtual reality, and metaverse applications. As the concept of a metaverse gains traction, avatars are poised to become even more ubiquitous; however, the public still faces barriers to creating personalized avatars, such as limited access to high-end capture hardware and/or the artist tools/expertise required to hand-craft such avatars.

With democratization in mind, we focus on the utilization of technologies that are currently widespread and are expected to be ubiquitous going forward, such as cell phones and webcams. Our pipeline intentionally avoids the use of detailed depth data and scanning, which has slowly been gaining popularity (e.g. because of the Kinect) but also at the same time losing ground in important hardware devices (e.g. the quality of depth data from the iPhone front-facing TrueDepth camera has been decreased in recent versions, and will likely be further decreased minimizing its use to security for the phone unlock feature). On the other hand, our pipeline embraces the continually improving RGB cameras in cell phones and webcams. We mostly rely on tasks executable by a non-expert user, e.g. asking for images from approximate angles, a turn-table video made while maintaining a neutral expression, and images of expressions that are prompted via a visual Simon Says approach.

Although there are a variety of recent and interesting approaches to representing facial avatars (e.g. using NeRFs [Mildenhall et al. 2021], PIFu [Saito et al. 2019], etc.), our pipeline reconstructs explicit geometry and de-lit texture so that it can be utilized in the widest number of existing graphics applications. The high-end special effects companies (e.g. ILM, Weta, Digital Domain, Pixar) typically use proprietary in-house software; however, there have been various efforts to democratize facial animation. Maya [Autodesk, INC. 2024] has various plugins that make facial animation easier, but one still has to sculpt/provide the face geometry/texture (and joints need to be placed by hand). Rigify [Blender 2023] in Blender [The Blender Foundation 2023] similarly, requires that one provide geometry/texture and place all the joints by hand. There are two standalone softwares, Headshot2 from Reallusion [Reallusion 2023] and FaceGen from Daz3D [Daz3D 2023], that create animatable geometry/texture without requiring the user to provide it, but the quality is hindered by a reconstruction that only uses a single front-facing image (no side or three-quarters images). MetaHumans [MHC 2023] in the Unreal Engine [Epic Games 2023] is the only publicly available tool that enables the creation of animatable geometry using more than just a single front-facing image, i.e. Mesh2Metahuman [M2M 2023] creates an animatable MetaHuman from input 3D geometry; unfortunately, there is no algorithm for obtaining texture (and one has to hand select from a library of defaults/sliders).

Although there is a plethora of prior work that creates floating face masks (often without texture, making it useless for preserving likeness), this geometry needs to be imported into animation software for subsequent use. Missing features (e.g. eyes, ears, inner lip and mouth, hair, etc.) need to be added manually or hallucinated, and the geometry needs to be retopologized (which can be lossy). Since this hallucination/lossiness modifies the likeness, it is essential that one be able to improve upon the likeness of the facial model after any such import. Methods that infer geometry/texture from photos cannot do this unless they are trained to infer animatable models directly (see e.g. [Shi et al. 2020][Lin et al. 2021]). Methods that solve an optimization problem to determine vertex positions (or parameterizations, e.g. morphable models [Blanz and Vetter 1999], FLAME[Li et al. 2017], etc.) and texel colors can typically be modified to improve the likeness of the facial model after the import.

It is well-known that geometry/texture reconstruction is best evaluated from novel views (see figure 1), but less discussion occurs around the fact that animated expressions further highlight inaccuracies in geometry reconstruction and texture alignment (see figure 2). A more subtle but also problematic issue is the mismatch between the template expressions on the animation rig and the corresponding expressions of an individual. The way someone smiles,

2 • Zhu, Y. et al

frowns, speaks, etc. is part of their motion signature and intrinsic to their likeness. We address this by using a Simon Says approach to prompt the user into producing various expressions that are reconstructed with our warping/projection approach; subsequently, the results are used to create a person-specific animation rig with motion signatures that better preserve the likeness of an individual.



Fig. 1. From left to right: geometry, same geometry textured from a frontfacing image, same geometry/texture as seen from a three-quarters view. This emphasizes how misleading a textured geometry can be when not considered from significantly novel views.



Fig. 2. From top left to top right: avatar geometry (derived from the geometry in Figure 1), same geometry textured from a front-facing image, same geometry/texture with a smile expression. This emphasizes how misleading a textured geometry can be when not considering various expressions. Bottom: zoomed-in view of the top middle and top right figures. Note in particular how part of the bottom lip texture (and the crease between the lips) appears on the top lip.

The novel contributions of our pipeline can be summarized as follows:

- We intentionally use baked-in lighting to create surrogate features by warping and projecting real-world images into the synthetic geometry's texture before subsequent optimization (of the geometry); notably, each view gets its own surrogate features (i.e. texture).
- We import our reconstructed geometry into the MetaHuman animation system, noting that this import is both lossy (i.e. modifies the input geometry) and hallucinates various

features (e.g. eyes, ears, inner lip and mouth, hair, etc.). Notably, this is the only publicly available option (we are aware of) that non-experts have for creating an animatable avatar from an input mesh. Subsequently, we show that our warping/projection approach can be used to improve the likeness of this imported animatable geometry.

- We propose a method for creating well-aligned de-lit textures from warped and projected images, preserving details such as moles, freckles, and stubble (important for maintaining likeness) while removing baked-in lighting. Importantly, our approach can be used to provide textures for animatable MetaHumans (which currently only utilize a library of defaults/sliders).
- We morph a default animation rig (similar to [Lin et al. 2022]) to match the neutral expression geometry, and subsequently use a visual Simon Says approach to prompt the user into producing various expressions that are reconstructed (via our warping/projection approach) and used to build a person-specific animation rig (preserving motion signature likeness).

Section 3 discusses how we obtain an initial reconstruction of the geometry in order to bootstrap our pipeline. Section 4 and Section 5 discuss our novel techniques for creating high-efficacy neutral/expressionless geometry and high resolution de-lit textures respectively. Notably, the result of Section 4 and Section 5 is a fully animatable rig (with eyes, mouth, hair, etc.), not just a floating face mask. Section 6 addresses the use of our method when one does not have access to the subject that one wants to reconstruct. Section 7 discusses the fitting of a template animation rig to the reconstructed geometry, as well as how our Simon Says framework can be used to create a person-specific animation rig that better preserves an individual's motion signatures.

2 RELATED WORK

In addition to the works discussed below, we refer the reader to the following survey papers: [Chrysos et al. 2018] [Zollhöfer et al. 2018] [Tewari et al. 2022] [Tretschk et al. 2023].

High-End Light Stages: The light stage [Debevec et al. 2000] [Debevec 2012] [Liu et al. 2022] uses hundreds of light/camera combinations to acquire a 4D reflectance field, enabling the highest quality facial reconstruction currently available. The light stage and similar systems [Ghosh et al. 2011] [Joo et al. 2017] [Hendler et al. 2018] [Zhang et al. 2022c] [Bolkart et al. 2023] that directly control illumination capture extremely high-fidelity geometry and separate albedo/specular/displacement maps. In addition to capturing a neutral/expressionless geometry and texture, these systems can be used to capture various facial expressions; subsequently, a visual effects artist (after some manual cleanup) can construct a high-quality facial animation rig. Moving towards democratization, various efforts have focused on reducing the required number of lights and/or cameras while maintaining high-quality results (e.g. [Zhang et al. 2004]); in particular, [Lattas et al. 2022] uses (easily portable) commodity components to construct facial capture systems that sit on a typical office desk.

High-End Multi-view: Various high-end systems [Beeler et al. 2022] [Beeler et al. 2010] [Bradley et al. 2010] [Riviere et al. 2020] lack the high degree of lighting control available to a light stage, but still use a number of carefully calibrated cameras (intrinsics and extrinsics) along with measured/controlled lighting in a laboratory setting. Such systems are often preferable (to light stages) for performance capture, since a light stage can feel constraining to an actor.

Single-view Democratization: The typical scenario for democratization is one in which the user has access to a single camera or cellphone (perhaps with a tripod) and no control over the lighting (but perhaps access to a cheap version of a chrome sphere, such as an ornament); interestingly, it is possible to program an at-home large-screen TV in order to mimic various light patterns from a light stage (see also [Sengupta et al. 2021]). In this democratized setting, noise and other inaccuracies (imperfect camera calibration, unknown lighting, etc.) can lead to disappointing results; thus, a strong prior on allowable face shapes is typically required. Although the 3D Morphable Model (3DMM) [Blanz and Vetter 1999] has historically been the most prevalent method of choice, recent works (such as FLAME [Li et al. 2017] and [Booth et al. 2016] [Tewari et al. 2017] [Tran and Liu 2018] [Wang et al. 2022] [Chandran et al. 2022]) have aimed to improve upon 3DMMs. Either optimization [Romdhani and Vetter 2003] [Li et al. 2017] or deep learning [Richardson et al. 2016] [Zhu et al. 2016] [Bulat and Tzimiropoulos 2017] [Dou et al. 2017] [Tewari et al. 2018] or both [Sanyal et al. 2019] can be used to regress the parameters of such models in order to best match an image. See also [Jackson et al. 2017] [Zeng et al. 2019] [Sengupta et al. 2018]. In addition, some works (e.g. [Richardson et al. 2016] [Tewari et al. 2018] [Tran et al. 2019] [Dib et al. 2023]) capture a residual geometry/texture on top of the "best-fit" model parameters.

Multi-view Democratization: The more successful multi-view techniques are typically pursued in the context of high-end systems (perhaps, most notably [Beeler et al. 2022]), since multiple cameras (often with carefully calibrated extrinsics) are required in order to obtain high-quality results. It is much more difficult to work with uncalibrated camera extrinsics and a single camera (with multiple images or video input). Most of these methods obtain no more than a simple floating face mask (sometimes with eyes), see e.g. the photometric stereo approaches in [Kemelmacher-Shlizerman and Seitz 2011][Kemelmacher-Shlizerman 2013][Liang et al. 2016], the structure from motion approaches in [Garg et al. 2013][Shi et al. 2014] [Ichim et al. 2015], the optimization approaches in [Blanz and Vetter 2003][Garrido et al. 2016][Piotraschke and Blanz 2016][Roth et al. 2016], and the neural network approaches in [Dou and Kakadiaris 2018] [Tewari et al. 2019] [Wu et al. 2019] [Bai et al. 2020] [Chaudhuri et al. 2020]. Of these, only [Tewari et al. 2019][Chaudhuri et al. 2020] aim to generate de-lit textures (the others either use no texture or simply splat the image onto the geometry). Only [Ichim et al. 2015] generates more than a floating face (they generate a full head, eyes, mouth, and animation rig, although the hair and eyebrows are incorrectly flattened into the texture). Most of these works focus on frontal or near-frontal views (occasionally showing obviously distorted side views); in fact, only [Liang et al. 2016][Tewari et al. 2019] show comparisons with ground truth side and/or three-quarters views.

Texture: In addition to geometry, texture is also acquired in many of the aforementioned approaches; however, there are a number of works that focus primarily on texture acquisition. [Kim et al. 2021] generate textures with baked-in lighting from a single image by training machine learning models in an unsupervised manner, and [Slossberg et al. 2022] adopts a similar method while removing the baked-in lighting. [Smith et al. 2020] [Han et al. 2023] aim to build a morphable face albedo/reflectance model by leveraging high quality texture data from high-end capture systems, and [Feng et al. 2022] [Ren et al. 2023][Rainer et al. 2023] train neural networks to reconstruct textures.

Building Animation Rigs: High-end special effects companies have proprietary in-house software, and democratized efforts are primarily limited to various Maya [Autodesk, INC. 2024] plugins, Rigify [Blender 2023] in Blender [The Blender Foundation 2023], Headshot2 from Reallusion [Reallusion 2023], FaceGen from Daz3D [Daz3D 2023], and MetaHumans [MHC 2023] in the Unreal Engine [Epic Games 2023]; in particular, Mesh2Metahuman [M2M 2023] is the only publicly available tool that enables the creation of animatable geometry using more than just a single front-facing image. It is also worth noting [Shi et al. 2020] [Lin et al. 2021], which train a network to infer animatable models directly. In addition to the identity priors used to regularize neutral/expressionless geometry capture (discussed above), both 3DMM [Blanz and Vetter 1999] and FLAME [Li et al. 2017] (and other models) have a separate set of blendshapes meant for facial animation (see [Lewis et al. 2014] for a review). Most of the aforementioned prior works that build animation rigs use these more academic models, perhaps since they are only aiming for technical demonstrations and not (actual) democratization. In real-time applications, joint transforms are often used instead of or in conjunction with blendshapes [Ward 2004]. The neutral expression is first deformed with joint-based linear-blend skinning, and blendshapes are (optionally) subsequently used as correctives (see e.g. [Franc 2023]).

Neural Rendering and Implicit Representations: Although we aim to build explicit geometry and de-lit textures so that they can be widely used in a variety of existing graphics pipelines, recent and interesting results in avatar generation have used NeRFs [Mildenhall et al. 2021] and implicit methods (such as [Saito et al. 2019] [Yariv et al. 2020] [Li et al. 2022] [Alldieck et al. 2022]) to represent 3D models. See e.g. [Sevastopolsky et al. 2020] [Zhang et al. 2022b] [Gafni et al. 2021] [Guo et al. 2021] [Wang et al. 2021] [Cao et al. 2022] [Gao et al. 2022] [Zhang et al. 2022a] [Zheng et al. 2023a] [Lin et al. 2023]. Notably, implicit representations can be converted to explicit representations (see e.g. [Azinović et al. 2023][Zheng et al. 2023b][Wang et al. 2023]), although the resulting explicit representations have not yet demonstrated the efficacy required in order to be widely adopted; of particular interest, [Wang et al. 2023] is motivated by light stage democratization, but their reconstructed geometry is difficult to evaluate since only frontal facing results are shown. There have also been attempts at building neural animation rigs, see e.g. [Qin et al. 2023].

3 INITIAL RECONSTRUCTION

We utilize two separate methods for the initial reconstruction of the geometry in order to bootstrap our process. The first method (presented in this section) assumes that one has access to a modern cellphone, and requires (a non-expert user) taking a few pictures from different angles and distances. The second method (see section 6) is more appropriate when one lacks access to the subject (e.g. one might desire a younger version of themself, the subject may be deceased, etc.), and only requires a short video of the subject (e.g. from a webcam, YouTube, etc.).



Fig. 3. From left to right: (a) captured image with a (blue) tracing of the silhouette, nostril, and lip corner, (b) initial triangle mesh created from the front view, (c) pixel-aligned projection of the front view triangle mesh onto the rough scan (with the aid of Laplacian smoothing), (d) accounting for silhouette boundaries of adjacent views, (e) MetaHuman reconstruction (note how fitting to a template hallucinates and modifies geometry).

We chose an Apple iPhone 13 Pro in order to demonstrate the process (although it is straight forward to extend these techniques to other phones). The back-facing dual camera was used to capture stereo color images from five views (front, left/right three-quarters, and left/right profile). For each view, stereo block matching [Konolige 1998] was used to obtain a crude depth estimate; then, a pixelaligned signed distance function was computed on a voxelized view frustum [Rasmussen et al. 2003], the voxelized view frustum was resampled onto a Cartesian grid, the fast marching method [Sethian 1999] was used to obtain signed Euclidean distance on the Cartesian grid, marching tetrahedra [Doi and Koide 1991] was used to construct a triangulated surface mesh, and segmentation [Yu et al. 2018] on the color images was used to remove background vertices (not corresponding to the subject). See Figure 3(b). To rigidly align the five triangle meshes, landmark detection [Hyprsense 2021] was run on one image from each view and the point on the triangle mesh corresponding to each visible landmark was identified (via ray tracing); then, the Procrustes algorithm [Gower 1975] was used to coarsely align view-adjacent pairs of meshes and iterative closest point [Holz et al. 2015] was used to refine this coarse alignment.

In order to refine the five triangle meshes and combine them into a single mesh, we obtained a rough scan from the iPhone front TrueDepth camera. For each view, one of the camera/image pairs was used to perturb the triangle mesh vertices in the pixel-aligned look-direction in order to best match the surface of the rough scan. These perturbations were regularized using Laplacian smoothing [Field 1988] on the displacements; importantly, this smoothing helps to alleviate mismatches between the triangle mesh geometry and the rough scan (e.g. mismatching nostril vertices with the cheek portion of the scan) while also providing displacements along silhouette boundaries that do not overlap with the scan. See Figure 3(c). Afterwards, each mesh was again perturbed in the pixel-aligned look-direction (again with Laplacian smoothing) in order to match the silhouette boundaries of its (one or two) adjacent view(s). See Figure 3(d). Finally, screened Poisson surface reconstruction [Kazhdan and Hoppe 2013] was used on the point cloud formed by the combined vertices of all five triangle meshes.

Along the lines of [Lin et al. 2022], we photon map the mesh to obtain a texture suitable for the Mesh2Metahuman pipeline [M2M 2023]. The Mesh2Metahuman pipeline retopoligizes the mesh to be consistent with an underlying MetaHuman animation rig [MHC 2023] (see section 7.1). Since the resulting neutral identity blend-shape $N^{\rm mh}$ may be far from the input mesh (both due to a lack of variety of scanned identities in the dataset and regularization), the Mesh2MetaHuman pipeline also outputs displacements $D^{\rm mh}$ that perturb the vertices so that $N^{\rm mh} + D^{\rm mh}$ is closer to the input mesh. Of course, one could use a different dataset and/or different optimization scheme; in fact, we found that subsequently optimizing to minimize the displacements (with a per-vertex loss) produced a neutral identity blendshape N significantly closer to $N^{\rm mh} + D^{\rm mh}$ than $N^{\rm mh}$ (resulting in an N + D with a smaller D).

It is worth noting that a number of prior works can achieve results similar to that shown in Figure 3(d) and thus obtain results that are similar to Figure 3(e) via the Mesh2Metahuman pipeline (perhaps using [Lin et al. 2022] or similar methods to obtain the input texture required to execute the Mesh2Metahuman pipeline). However, regularized and hallucinated geometry will adversely affect the likeness, the Mesh2Metahuman tool does not contain a method for obtaining a de-lit texture, and the default animation rig will not actuate in accordance with the motion signatures of the subject. We address these issues in section 4, 5, and 7 respectively.

4 GEOMETRY REFINEMENT

In a democratized pipeline, camera intrinsics might be available (to some accuracy); however, knowledge of lighting, albedo, and other information required for disentanglement of the geometry from the texture and lighting will be lacking. In particular, a synthetic rendering of the current guess for the geometry will have different features than the real-world image because of mismatches in geometry, texture, and lighting. Our key insight is that one can bake (entangled) lighting and texture information from the real-world image onto the current guess for the geometry in order to provide surrogate landmark information (similar in spirit to painting on mo-cap dots) for subsequent geometry optimization.

The main issue with baking in lighting and texture in order to create surrogate features is that there is a mismatch between the synthetic and real-world geometry. Thus, after an initial rigid alignment, each real-world image is non-linearly warped to better align with the synthetic geometry before projecting the pixel colors from the real-world image into the texture of the synthetic geometry. See Figure 4. Afterwards, the synthetic geometry is optimized to match the original unwarped image with the aid of the surrogate features that have been baked into its texture. Although using multiple views and the appropriate (different) projected texture for each view is essential to the efficacy of this process, we initially describe the method in terms of a single view (in sections 4.1-4.3) before describing the modifications required to accommodate a multi-view approach (in section 4.4).

4.1 Rigid Alignment and Lighting

Given a current geometry and texture, we use landmarks [Hyprsense 2021] to rigidly align it with a real-world image; subsequently, a spherical harmonics lighting model [Ramamoorthi and Hanrahan 2001] is optimized to best match a synthetic rendering (of the geometry/texture) to the real-world image. This rigid alignment and estimated lighting is required in order to obtain commensurate results when using a segmentation map network [Yu et al. 2018] on both the synthetic rendering and the real-world image (segmentation maps are used in the construction of the non-linear warp).

Given a landmark tracker of choice, we hand-label a (non-sliding) subset of the landmarks on the MetaHuman template geometry (this hand-labeling only needs to be done once per choice of landmark tracker); then, a simple least-squares fit can be used to provide an initial estimate for the rigid alignment between any geometry and any image. Similar to [Wu et al. 2023a], we refine the initial estimate by comparing landmarks computed on renderings of the geometry as opposed to using the hand-labeled landmarks; importantly, this allows landmarks to slide and alleviates issues caused by consistent errors in the landmark tracker, such as always labeling a specific marker too high/low/etc. on similar images (these errors are more prevalent for profile angles and out-of-distribution images). Assuming the intrinsic camera parameters are known, we solve for the extrinsic parameters T (rotation and translation) to refine the rigid alignment by minimizing

$$\mathcal{L}_{rigid}(\mathbf{T}) = \sum_{l} \left[\zeta(\Psi(\mathbf{v}, \mathcal{T}, L^{a}, \mathbf{T}))_{l} - \zeta(\mathcal{I}^{R})_{l} \right]^{2}$$

where Ψ differentiably rasterizes (we use [Ravi et al. 2020]) the geometry (vertex positions v, texture \mathcal{T}) into an image using (only) ambient lighting L^a . ζ is a landmark tracker that operates on images, l are the computed landmarks ([Hyprsense 2021] computes 159 landmarks), and \mathcal{I}^R is the real-world image. Note that [Hyprsense 2021] defines each landmark as a weighted sum of heat-map values in order to preserve differentiability.

Although we found ambient lighting to be sufficient for the landmark tracker, it was insufficient for the segmentation map network from [Yu et al. 2018]. Therefore, we next estimate the parameters of a spherical harmonics lighting model L (while keeping T fixed) by minimizing a per-pixel objective function

$$\mathcal{L}_{light}(L) = \sum_{i \notin \mathcal{B}} \left[\Psi(\mathbf{v}, \mathcal{T}, L, \mathbf{T})_i - I_i^R \right]^2$$

that ignores background pixels.

4.2 Warping and Projection

A synthetic rendering of the current geometry and texture (using the spherical harmonics lighting estimate) will typically be quite different than the real-world image due to the so-called "domain gap"; moreover, it is difficult (if not impossible) to disentangle the errors in geometry, texture, and lighting from each other. Thus, we use semantic segmentation [Yu et al. 2018] of key coarse regions of the face (e.g. nose, lips, eyebrows, silhouette) as the signal for aligning a real-world image with a synthetic rendering. The semantic segmentation also helps to account for subtle differences in expression (alleviating the difficulties associated with capturing expressions with a generic model).

We construct an optical flow field **f** to smoothly warp a one-hot encoding of the semantic segmentation of the real-world image $\mathbb{1}(Seg(I^R))$ to match (as well as possible) a one-hot encoding of the semantic segmentation of the synthetic rendering $\mathbb{1}(Seg(I^S))$. Note that we merge the hair and body regions into the background category of the one-hot encoding in order to focus on the facial shape. The per-pixel objective function

$$\mathcal{L}_{seg}(\mathbf{f}) = \sum_{i} \left(\mathbb{1}(Seg(\mathcal{I}^{S}))_{i} - \eta(\mathbb{1}(Seg(\mathcal{I}^{R})), \mathbf{f})_{i} \right)^{2}$$
(1)

compares $\mathbb{1}(Seg(\mathcal{I}^S))$ to a non-linear warp of $\mathbb{1}(Seg(\mathcal{I}^R))$ where η warps one image to another using the flow field **f** and bilinear interpolation. Since \mathcal{L}_{seg} only provides penalties near differences in the segmentations, we add regularization to minimize the per-pixel norm of **f** and the per-pixel Laplacian (using a 5x5 stencil) of **f**.

Although using \mathcal{L}_{seg} and the regularizers behaved as expected, we found that post-processing f led to significantly improved results. The post-process solves a 2D Laplace equation (using a standard 5-point stencil) using Dirichlet boundary conditions on pixels with differently-labeled neighbors (according to $Seg(\mathcal{I}^R)$) and Neumann boundary conditions on the image boundary. This preserves **f** on the boundaries of $Seg(\mathcal{I}^R)$ while guaranteeing that it is smooth elsewhere.

Lastly, we project the warped (and thus better aligned) real-world image \mathcal{I}^W into the texture for the synthetic geometry using the photon mapping technique proposed in [Lin et al. 2022]. Each pixel of \mathcal{I}^W is projected onto the geometry to determine texture coordinates for storing a "photon" based on the pixel color. In the gathering step, the color for each pixel in the texture map is determined using a weighted average of the *k* nearest photons. Typically, the weights would decay with distance; however, we additionally scale the weights with a power law for specular falloff $(\max(0, -\mathbf{n} \cdot \mathbf{r}))^P$ where **n** is the (unit) normal to the geometry and **r** is the (unit) ray direction from the pixel to the geometry. This additional scaling diminishes the contribution from photons near occlusion boundaries, which may still be misaligned after the warp.

4.3 Optimization

Since the texture \mathcal{T} computed via section 4.2 already contains bakedin lighting, only ambient lighting L^a is required for the differentiable rasterizer Ψ when using inverse rendering to optimize the geometry. Perturbations of the (geometry) vertex positions v can be computed by minimizing a per-pixel objective function

$$\mathcal{L}_{photo}(\mathbf{v}) = \sum_{i \notin \mathcal{B}} \left[\Psi(\mathbf{v}, \mathcal{T}, L^a, \mathbf{T})_i - \mathcal{I}_i^R \right]^2$$

that ignores background pixels (note that the hair and body regions are still kept merged into the background).



Fig. 4. A real-world image (shown in the first figure) is warped (in image space) to better match a synthetic rendering of a current guess for the geometry (using an appropriate texture). A zoomed-in view of the real-world image is shown before (second figure) and after (third figure) the warp. The fourth figure shows the current geometry, and the fifth figure shows the result obtained by projecting the warped real-world image onto that geometry. This new texture contains baked-in lighting that provides surrogate features useful when optimizing the synthetic geometry to match the original unwarped image.

In order to encourage the matching of semantic information about facial shape, we additionally utilize a second objective function based on the semantic segmentations of \mathcal{I}^R and the synthetically rendered \mathcal{I}^S . Let \mathbf{v}^k be the geometry at the beginning of the kth iteration and $\mathcal{I}^{S,k}$ be the synthetic rendering of that geometry; then, the procedure from section 4.2 can be used to construct an optical flow \mathbf{f}^k that smoothly warps $\mathbbm{1}(Seg(\mathcal{I}^R))$ to $\mathbbm{1}(Seg(\mathcal{I}^{S,k}))$. Given sample points \mathbf{v}^k_j on the geometry \mathbf{v}^k , a ray tracer can be used to compute their screen-space locations $\tilde{\Psi}(\mathbf{v}^k_j, \mathbf{T})$; then, the per-geometry-sample objective function

$$\mathcal{L}_{sem}^{k}(\mathbf{v}) = \sum_{j} \left[\tilde{\Psi}(\mathbf{v}_{j}, \mathbf{T}) - \tilde{\Psi}(\mathbf{v}_{j}^{k}, \mathbf{T}) + \mathbf{f}_{j}^{k} \right]^{2}$$

compares the screen-space motion of each sample point to the optical flow \mathbf{f}^k that smoothly warps $\mathbb{1}(Seg(\mathcal{I}^{S,k}))$ to $\mathbb{1}(Seg(\mathcal{I}^R))$ at that location (\mathbf{f}^k_j is the interpolation of \mathbf{f}^k to the screen-space location $\tilde{\Psi}(\mathbf{v}^k_j, \mathbf{T})$). Although there are many potential choices for the sample points (e.g. one could choose all visible triangle vertices), we use the sub-triangle locations that rasterize to the center of each non-background pixel (removing the need to interpolate \mathbf{f}^k). This makes the sample points \mathbf{v}_j depend on the triangle vertices via barycentric interpolation.

Since the lighting is baked into the surrogate texture, neither \mathcal{L}_{photo} nor \mathcal{L}_{sem}^k has or requires any dependence on the normal vectors of v; thus, geometric regularization is essential. We utilize two geometric regularizers: \mathcal{L}_{edge} and $\mathcal{L}_{laplace}$. \mathcal{L}_{edge} penalizes changes in the edge lengths, computed by comparing edge lengths in the current mesh to the lengths of the corresponding edges in the initial pre-optimization mesh. $\mathcal{L}_{laplace}$ penalizes the Laplacian (using a one-ring stencil) of the vertex displacements, computed by comparing vertex positions in the current mesh to their corresponding positions in the pre-optimization mesh; importantly, this promotes smooth displacements, but does not smooth out features in the initial pre-optimization mesh as the Laplacian of the vertex positions would.

4.4 Using Multiple Views

Since both rigid alignment and segmentation maps struggle with non-front-facing views (especially with non-photorealistic textures), we begin the multi-view process by warping and projecting the frontfacing real-world image into the texture (as per subsection 4.2); then, this photorealistic texture is used to compute rigid alignment and segmentation maps for each of the other views. Although one might blend textures from the various views together into a single texture, this incorrectly averages lighting information (which varies according to view); thus, we compute (as per subsection 4.2) and maintain separate textures for each view, and subsequently optimize one geometry with separate losses (\mathcal{L}_{photo} and \mathcal{L}_{sem}^{k} , each using the view-appropriate \mathcal{T}) for each view combined into the same objective function (noting that one could optionally increase the weights on \mathcal{L}_{edge} and $\mathcal{L}_{laplace}$ in order to balance the increased forces from using multiple views). In addition, the rigid alignment and segmentation maps seemed to perform best when using the texture from the previous iteration corresponding to the view under consideration (when that texture exists).

The entire multi-view approach can be repeated iteratively, obtaining a new geometry after each iteration. We obtained the best results using only a single front-facing view in the first iteration and multiple views in subsequent iterations. Since this iterative approach jointly optimizes both the geometry and the camera extrinsics (for the utilized real-world images), it can be beneficial to repeat the initial image-based reconstruction of the geometry (in section 3) using these improved camera extrinsics; notably, we found that this led to a significant improvement in both the initial image-based reconstruction as well as the subsequent multi-view iteration (which one might expect given the efficacy of end-to-end approaches). See Figure 5.

5 DE-LIGHTING TEXTURE

The real-world entanglement of geometry, texture, and lighting coupled with human perception makes it difficult to ascertain the quality of a reconstruction (this is especially misleading in prior works that compare an entangled reconstruction to a reference image); however, disentangled results are required in order to use a reconstruction in a standard graphics pipeline (especially when the geometry/texture needs to be re-lit). See Figure 6.

An initial guess for the texture can obtained from the warping and projection discussed in section 4, the photon mapping algorithm discussed in [Lin et al. 2022], or any other prior work; however,



Fig. 5. From left to right: (a) captured image with a (blue) tracing of the silhouette, nostril, lip corner, eye corner, and mouth corner, (b) MetaHuman reconstruction from section 3 Figure 3, (c) the results obtained by using our geometry refinement process (in section 4). In particular, note the improvements in the eye and mouth regions.



Fig. 6. It is often hard to distinguish the efficacy of a floating, grey-shaded face mask (first column). Adding the head, neck, eyes, etc. gives a better indication of an identity (second column). However, only after adding textures and hair does human perception cause the images to start migrating into and through the uncanny valley (third column). The top row shows a state-of-the-art geometry reconstruction result from the Mesh2MetaHuman toolset, and the bottom row shows our reconstructed geometry. The Mesh2MetaHuman reconstruction requires an iPhone 12's depth sensor, which is significantly better than the depth sensors in the iPhone 13 and newer models (the Mesh2MetaHuman reconstructed geometries, since the Mesh2MetaHuman toolset lacks the ability to reconstruct textures disentangled from lighting and geometry.

most methods result in misaligned textures due to imperfections in the predicted geometry. Thus in section 5.1, we discuss how to leverage the methods proposed in section 4 in order to obtain better texture alignment. In section 5.2, we illustrate how lighting estimates can be used to remove baked in lighting via inverse rendering. In section 5.3, we discuss the separation of an acquired texture into high-frequency and low-frequency components and the subsequent projection of the low-frequency component into a pre-computed PCA basis (further removing baked in lighting information).

5.1 Texture Alignment and Averaging

Importantly, popular face representation frameworks (we use the Metahuman [MHC 2023]) contain at least one default texture (often synthetic) that has been correctly aligned with the texture coordinates on the geometry; thus, we aim to align our subject-specific texture with that default texture. At the end of the geometry refinement (in section 4), each real-world image can be projected into the texture map (and gathered to pixels/texels) in order to obtain a corresponding texture. These projected textures will typically have stretching artifacts and misalignment in areas where the camera look direction is at a grazing angle to the geometry. Aiming to eliminate these artifacts, we combine all the projected textures into a single stitched texture by choosing the color at each texel (texture map pixel) from the texture map that had the most orthogonal geometry (as compared to the camera look direction) at that texel. Afterwards, the stitched texture is warped to align with the default Metahuman texture; then, all of the projected textures are warped to match the aligned stitched texture. See Figure 7.



Fig. 7. The first figure shows color-coding to indicate which texture (front view: white, three-quarters view: light red/blue, side view: dark red/blue) is contributing to the stitched texture (in the second figure). The stitched texture is warped to align with the default MetaHuman texture (shown in the third figure), and then each of the individual textures is warped to align with the stitched texture; afterwards, the textures are averaged together using the blending discussed in section 5.1. The final result is shown in the fourth figure.

The warping (in texture space) is accomplished by computing a flow field **f** along the lines of section 4.2 with a few modifications. The segmentation map network can be run on a rendered image (using the texture under consideration) and subsequently projected back onto the geometry and into texture space; alternatively, the segmentation map network can be run directly on the texture image (and works surprisingly well). Since each real-world image only fills in a portion of the full texture, we found it beneficial to include an additional term

$$\mathcal{L}_{lmk}(\mathbf{f}) = \sum_{l} \left(\zeta_{l}^{\mathrm{mh}} - \eta(\zeta(\mathcal{I}^{R}), \mathbf{f})_{l} \right)^{2}$$

that encourages landmark alignment. $\zeta_l^{\rm mh}$ are the fixed locations in texture space of hand-labeled landmarks on the MetaHuman template. The landmark locations $\zeta(I^R)$ can be computed on the real-world image I^R and subsequently projected into texture space

8 • Zhu, Y. et al

or be computed directly on the projected texture (similar to the segmentation map). η warps the landmarks using bilinear interpolation (to the pre-warped sub-pixel landmark locations) of the flow field **f**.

When warping each of the individual textures to align with the stitched texture, we utilize (dark) moles as additional landmarks (in the above equation). A square sample window is used to identify texels that are darker than the mean texel color in the window. The center of each connected component is identified as a mole whenever the connected component's area is larger than a threshold. Mole correspondence (between the stitched texture and the texture that is being warped to match it) is established by a greedy algorithm that finds the closest mole with similar area working outwards from the tip of the nose.

After warping, the textures are combined together using a specular falloff weighting along the lines of that discussed at the end of section 4.2. Notably, since the face was rotated in a static lighting environment, this averaging helps to remove baked in lighting information. We additionally reduce the per-texel specular falloff weighting of a texture that has an adjacent view that is more orthogonal to the geometry at that texel, since the adjacent view can likely provide a better approximation to the texture. For each view, orthogonality is measured at each texel as $-\mathbf{n} \cdot \mathbf{r}$; importantly, this orthogonality measure is labeled to be unusable whenever it is negative. Using this measure of orthogonality, the specular falloff weighting is further multiplied by the ratio of the orthogonality measures between the current texture and an adjacent texture whenever that ratio is less than one (when there are two choices for the adjacent texture, the one with the smaller ratio is used).

5.2 Leveraging Lighting Estimates

An approximation to the real-world lighting can be used to improve upon the result from section 5.1. This can be achieved numerically using inverse rendering (e.g. to estimate a spherical harmonics approximation [Ramamoorthi and Hanrahan 2001]) or captured directly (e.g. via a chrome sphere environment map). Some phone apps (e.g. HDReye [HDReye Technologies Inc. 2023]) use HDR capture and automatic stitching to create environment maps; in addition, consumer-level 360° cameras (e.g. Insta360 X3 [Insta360 2023] and Ricoh Theta [RICOH360 2023]) are becoming more prevalent. We have also experimented with programming an at-home large-screen TV in order to mimic various light patterns from a light stage (see also [Sengupta et al. 2021]).

Given any of the aforementioned lighting estimates, we can improve the texture from section 5.1 via inverse rendering. The objective function includes the per-pixel differences between the realworld image and the corresponding synthetic rendering for each of the five views, but these per-pixel differences are given specular falloff weightings (discussed at the end of section 4.2) in order to diminish contributions along the silhouette boundaries of each view. A Laplacian term is used on the texel deltas to ensure smooth changes in the texture. The spherical harmonics parameters can be jointly optimized, if desired.

5.3 Frequency Separation and PCA Projection

Whether a texture is acquired via the algorithms proposed in section 5.1 or 5.2 or the photon mapping process from section 4 (or other previous work), the result needs to be converted into a fully de-lit albedo texture (required in a standard graphics pipeline). At this point in our process, the majority of the baked lighting information is contained in the lower frequencies of the texture, while personalized features such as moles tend to be contained in the higher frequencies. Thus, we separate the texture into high-frequency and a low-frequency components by applying a Gaussian low-pass filter.

We remove the baked-in lighting from the low-frequency component via PCA projection. Starting with the MetaHuman database of albedo textures (from real-world scans), we perform whitening and subsequently calculate a PCA basis. Then, we optimize the coefficients of the first five bases in order to match the pixel color of the whitened low-frequency component using an L2 regularizer. We ignore eyebrow regions (since they are incorrectly baked into the texture), eye folds/sockets, and the inner nostril during PCA projection, but do reconstruct these regions. In order to capture some of the facial stubble, moles, and other details, the high-frequency component is added back to the non-ignored regions. See Figure 8.



Fig. 8. PCA-projected low frequency texture (left). Final texture (right), after adding back the high-frequency component. Bottom row is a zoomed-in version of the top row. Note the facial stubble, lip wrinkles, and mole.

We select the best-fit hair and eyebrows from the MetaHuman preset options. Since the MetaHuman database has relatively limited choices for hair and eyebrows (which one would expect to improve over time given the variety of character customization abilities in recent high-end video games), we slightly (manually) sculpt the hair and eyebrows to better match the images. The character customization controls for the eyes are more readily usable (no manual sculpting is necessary). Notably, character customization options are becoming more democratized, and non-expert users (non-artists) can reasonably match hair, eyebrows, etc. (when appropriate options are available); however, a non-expert user (or even a trained artist) would struggle to capture the overall face shape and appearance (which is the focus of Sections 4 and 5). See Figures 9, 10, and 11.

Democratizing the Creation of Animatable Facial Avatars • 9



Fig. 9. Comparison between captured images and our reconstructed geometry/texture for five views. The synthetic rendering utilizes lighting similar to that present in the real-world images. See also Figure 11.



Fig. 10. Comparison between captured images and our reconstructed geometry/texture for five views. The synthetic rendering utilizes lighting similar to that present in the real-world images. See also Figure 11.

6 VIDEO-BASED RECONSTRUCTION

In this section, we briefly discuss how our method adapts to the situation when one lacks access to the subject (e.g. one might desire a younger version of themself, the subject may be deceased, etc.). When high-quality images of a neutral expression in each of our five required views are available (either as images or as a subset of a webcam/YouTube/etc. video), the methods proposed in sections 4 and 5 require no modification; however, the initial reconstruction (from section 3) would need to be replaced with template MetaHuman geometry (typically degrading the quality of the final result). A better initial guess could be obtained by training a deepfake model on the video footage (as discussed in [Lin et al. 2022]). When images of a neutral expression in each of our five required views are not available, a deepfake model (trained on the video footage) could be used to create surrogate images suitable for our pipeline. Figure 12 illustrates the results we obtained on two publicly recognizable figures.

7 ANIMATION RIG

The Mesh2Metahuman pipeline creates geometry topoligized to be consistent with an underlying animation rig (see section 7.1). As noted near the end of Section 3, additional displacements (on top of



Fig. 11. The reconstructed avatars from Figures 9 and 10 rendered under various novel lighting conditions.

the neutral identity blendshape) are typically required in order to better match an input mesh. When these additional displacements are large (see e.g. Figure 14), subsequent animations often possess undesirable artifacts. Section 7.2 discusses how a volumetric morph of the animation rig removes the need for these additional displacements. Although the morphed animation rig fits the reconstructed geometry quite well, it still will not have the same motion signatures as the real-world subject (e.g. the subject's actual smile will be dissimilar to a dialed-in smile on the animation rig). We address this by capturing a number of basic expressions via a Simon Says approach (Section 7.3), and subsequently using them to build a personalized animation rig (Section 7.4).

7.1 MetaHuman Animation Rig

The MetaHuman animation rig [Franc 2023] is freely available in the Unreal Engine [Epic Games 2023], and the open-source code can be used outside the Unreal Engine as well. The top level implementation exposes a set of sliders that non-experts can readily use to control a face, see Figure 13. For the sake of real-time applications, deformations from a neutral geometry are primarily implemented via joints and linear blend skinning. The approximately 200 control sliders are mapped to translations and rotations of about 800 joints.



Fig. 12. In order to demonstrate the efficacy of our approach (from sections 4 and 5) when one does not have access to the subject, we chose two publicly recognizable figures to reconstruct from online images/video.

A more computationally expensive version of the animation rig uses about 700 blendshape correctives on top of the joint-based deformations.

Given a textured face geometry (the texture needs to be suitable for landmark detection), the Mesh2MetaHuman pipeline [M2M 2023] creates an appropriately topologized neutral identity blendshape $N^{\rm mh}$ as well as a corresponding animation rig. The placements of the joints, the mappings from slider controls to joint displacements (position and orientation), the linear blend skinning weights, the blendshape correctives, and the mappings from slider controls to blendshape corrective weights are all determined automatically by leveraging a database of expertly hand-crafted animation rigs. As is typical, this database was constructed from light stage (and other) scans of various individuals. Due to limited representability in the database, $N^{\rm mh}$ can deviate quite a bit from the input geometry; thus, an additional displacement blendshape D^{mh} is added on top of $N^{\rm mh}$. See Figure 14. Although $D^{\rm mh}$ mostly rectifies the inability of the database to adequately represent the input geometry, the animation rig is meant for $N^{\rm mh}$ not $N^{\rm mh} + D^{\rm mh}$. See Figure 15. This can lead to various animation artifacts (we remedy this in section 7.2). See Figure 16. It is also important to note that one would not



Fig. 13. 2D graphical user interface for MetaHuman animation rig control.

generally expect the blending of animation rigs from a database to adequately capture the motion signatures of any particular individual not already in that database (we address this in Sections 7.3 and 7.4).

7.2 Volumetric Morph

Figure 14 shows input geometry that is not well-represented by the MetaHuman database. This leads to a reconstructed neutral blend-shape $N^{\rm mh}$ that differs significantly from the input geometry; thus, additional displacements $D^{\rm mh}$ are added on top of $N^{\rm mh}$. Problematically, the animation rig is intended for $N^{\rm mh}$ not $N^{\rm mh} + D^{\rm mh}$. See Figures 15 and 16. We remedy this by volumetrically morphing the animation rig to better fit $N^{\rm mh} + D^{\rm mh}$, allowing for the removal of the extra displacements entirely. In particular, we morph from $N^{\rm mh}$ to \hat{N} , where \hat{N} represents the result of our geometry optimization (from Section 4).

Following [Cong et al. 2015], we extend the $N^{\rm mh}$ to \hat{N} per-vertex displacements to a volumetric field by solving three decoupled threedimensional Poisson equations. The computational domain is specified by an oversized bounding box and discretized with a Cartesian grid of suitable resolution. Dirichlet boundary conditions are specified on any Cartesian grid edge that intersects $N^{\rm mh}$ using barycentric interpolation of the per-vertex displacements. After solving the Poisson equations, trilinear interpolation can be used to determine the displacements required to morph each joint's position. In order to update a joint's orientation, three additional sample points ($x_0 + \epsilon i$, $x_0 + \epsilon_j$, and $x_0 + \epsilon_k$ where x_0 is the unmorphed joint position and ϵ is a small number) are morphed to their new locations and subsequently used to determine a rotation (taking care to re-orthogonalize after the morph). The blendshape correctives can be rewritten in the local coordinates of N^{mh}, identity-transformed to the local coordinates of \hat{N} , and then transformed back out of local coordinates (independent of the volumetric morph). The mappings from slider controls to

Democratizing the Creation of Animatable Facial Avatars • 11



Fig. 14. Top: scan of an individual that is not well-represented by the MetaHuman database. Bottom left: identity blendshape $N^{\rm mh}$ reconstructed from the MetaHuman database. Bottom right: adding an additional displacement blendshape $D^{\rm mh}$ to the identity blendshape $N^{\rm mh}$ results in a better match to the scan.

joint displacements (position and orientation), the linear blend skinning weights, and the mappings from slider controls to blendshape corrective weights are left unchanged.

7.3 Simon Says

One would not generally expect the blending of animation rigs from a database to adequately capture the motion signatures of any particular individual not already in that database; thus, we capture a number of basic expressions and subsequently use them to build a personalized animation rig (see Section 7.4). To balance between animation rig quality and capture session length (and complexity), we select 27 expressions (divided into two groups):

- pucker, nose wrinkle, cheek raise, mouth stretch, squint lower eyelid, lip corner pull, jaw open, brow lower, brow raise, blink
- upper lip raise, nose wrinkle combined with upper lip raise, sharp lip corner pull, mouth dimple, lip corner depress, lower lip depress, purse, lips towards, funnel, funnel purse, funnel towards, oh, jaw open extreme, lip corner pull combined



Fig. 15. Top: the Mesh2MetaHuman animation rig fits $N^{\rm mh}$ well, by design. Middle: the Mesh2MetaHuman animation rig does not fit $N^{\rm mh} + D^{\rm mh}$ well. Joints can be too deep inside of or even outside of the surface geometry; perhaps more importantly, joints are improperly aligned with the surface topology. Bottom: our volumetrically morphed animation rig fits $N^{\rm mh} + D^{\rm mh}$ well.

with jaw open, mouth stretch combined with jaw open, smile, smile stretch combined with jaw open



Fig. 16. Top row: the animation rig fits $N^{\rm mh}$ well, and also animates it well. Middle row: The animation rig does not fit $N^{\rm mh} + D^{\rm mh}$ well, and thus does not animate it well. Bottom row: our morphed animation rig fits $N^{\rm mh} + D^{\rm mh}$ well, and also animates it well. All three rows have identical animation rig controls (for purse and funnel) dialed in. Note the semantically similar expressions in the top and bottom rows (as expected); in contrast, the poor-fitting animation rig (middle row) incorrectly makes a semantically dissimilar expression.

noting that the second group can (optionally) be omitted for brevity.

The capture session is guided via a Simon Says approach, where the user's avatar (built via Sections 3, 4, 5, and 7.2) makes an expression that the user attempts to match. See Figure 17. Instead of giving users obscure terminology for various expressions, visual cues better enable non-experts to quickly understand each expression (especially when the visual cues are on an avatar similar in appearance to the user). The visual interface shows a video feed of the user side by side with their avatar, see Figure 18.



Fig. 17. The 27 expressions we capture in order to build a personalized animation rig. Left: the avatar makes an expression, which is shown to the user (see Figure 18). Right: the user makes their version of that expression. The differences between the animation rig expressions and the user's expressions are quite apparent; in particular, the user opens their jaw differently, makes different mouth shapes, and sometimes struggles to make a similar expression (for example, see the last column of the first row).



Fig. 18. Visual interface for the Simon Says capture session. A video feed of the user is shown side by side with their avatar. When a pose is finalized, the user clicks a button in order to capture/save the image of their face.

7.4 Personalized Animation Rigs

For each of the 27 images from the Simon Says capture session (see Figure 17, columns 2,4, and 6), the geometry refinement method from Section 4 can be used to reconstruct geometry corresponding to the expression. Since only a single image (not multi-view) is available, more regularization is required; thus, we optimize for animation rig degrees of freedom instead of per-vertex displacements (but otherwise follow Section 4). Using animation rig degrees of freedom also allows one to ignore or remove spurious geometry matching unrelated to a given expression (for example, only mouth degrees of freedom should be activated during a smile). Degrees of freedom that should not be active for a given expression can either be ignored during the optimization or be set identically to zero as a post process (after the optimization). In addition, this allows the capture of symmetric expressions (such as a smile) to be used to obtain geometry for asymmetric expressions (such as a left-only or right-only smile) simply by setting (approximately) half of the controls to zero as a post process.

For each of the 27 basic expressions, the reconstructed geometry is used to modify the corresponding expression in the animation rig. Since each of these is a primary expression in the animation rig, they are straightforward to modify. After dialing the sliders corresponding to an expression to their maximum value, the difference between the deformed avatar and the reconstructed geometry is used to replace the corrective blendshape associated with that expression. Optionally, one could instead modify the joint displacements in order to be more consistent with the reconstructed geometry; although, a (smaller magnitude) corrective blendshape would still be required in order to exactly match the reconstructed geometry. Modifying primary expressions also indirectly modifies the complex expressions that are composed of multiple primary expressions. The complex expressions can be preserved by modifying their corrective blendshapes with displacements opposite of those that were used to modify any primary expression that they depend on. A similar approach could be taken when choosing to modify joint displacements. See Figure 19for two samples (chosen from the 27, for brevity).



Fig. 19. Before (left) and after (right) modifying the animation rig with reconstructed geometry from the Simon Says capture session (the corresponding image from the Simon Says capture session is shown in the middle).

8 DISCUSSION AND CONCLUSIONS

We chose to use the MetaHuman framework because it provides a full head, face, and neck (not just a floating face mask) as well as a default animation rig (however, we did have to devise our own method for adding textures). The use of any other suitable database (either existing or future) will be subject to issues similar to those discussed in this paper: the size of the database will be limited by the difficulties associated with scanning individuals and creating suitable personalized animation rigs, the representability of the database will be lacking due to its limited size (especially when considering the large variance in human faces and expressions), the limited representability will lead to poor results (poor geometry, texture, and animation rigs) for individuals that are not well-represented by the database, etc. Our methods for geometry refinement (Section 4), creating de-lit textures (Section 5), morphing template animation rigs to better fit optimized geometry (Section 7.2), and modifying animation rigs to contain personalized motion signatures (Sections 7.3 and 7.4) should thus be useful for improving the results from any database (not only the MetaHuman database)

We would like to stress the importance of capturing and utilizing profile views, as they are crucial for obtaining the correct geometry (not discernible from front facing views, due to depth ambiguities). Although there is a plethora of work that focuses on primarily front facing views, the importance of capturing and evaluating the results from novel views is becoming more prevalent in the literature, e.g. [Wu et al. 2023b]. The importance of obtaining the correct geometry (and texture) becomes even more apparent when subsequently animating/deforming the face; although the neutral identity may appear correct, expressions can lie in the uncanny valley. This is perhaps even more important when one considers subsequent biomechanical simulations where areas and volumes (not only two-dimensional projections) need to be accurate.

Our approach of using baked-in lighting to create surrogate features can be used to improve an initial reconstruction from any method. More importantly, it can be used to improve avatars (such as MetaHumans) directly. This is essential since any mesh-to-rig pipeline will be subject to the inadequacies of its database (and thus both hallucinate and be lossy). Our method for obtaining de-lit textures allows the reconstructed avatar to appear realistic in novel lighting environments while still maintaining the high-frequency moles, stubble, etc. required to preserve likeness. Our video-based reconstruction results illustrate that we can create avatars for subjects we do not have access to with only a small modification to our pipeline. Finally, our rig building approach alleviates some of the issues associated with the representability gap between any particular individual and a database of (preferably, carefully handcrafted) rigged avatars. The volumetric morph better fits the animation rig to the geometry (removing both deformation artifacts and issues with semantics), and the Simon Says approach enables a non-expert user to capture expressions important for reproducing their motion signatures.

9 ACKNOWLEDGEMENTS

Research supported in part by ONR N00014-19-1-2285, ONR N00014-21-1-2771, and Epic Games. We would like to thank Reza and Behzad at ONR for supporting our efforts into machine learning. We would also like to thank Dragan Davidovic, Iain Matthews, Jihun Yu, Jovan Mijatov, Kamy Leach, Kim Libreri, Michael Lentine, Relja Ljubobratović, Steven Caulkin, Thibaut Weise, Jungwoon Park, and Vladimir Mastilovic for their insightful discussions.

REFERENCES

- Thiemo Alldieck, Mihai Zanfir, and Cristian Sminchisescu. 2022. Photorealistic monocular 3d reconstruction of humans wearing clothing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1506–1515.
- Autodesk, INC. 2024. Maya. https:/autodesk.com/maya
- Dejan Azinović, Olivier Maury, Christophe Hery, Matthias Nießner, and Justus Thies. 2023. High-res facial appearance capture from polarized smartphone images. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 16836–16846.
- Ziqian Bai, Zhaopeng Cui, Jamal Ahmed Rahim, Xiaoming Liu, and Ping Tan. 2020. Deep Facial Non-Rigid Multi-View Stereo. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 5850–5860.
- Thabo Beeler, Bernd Bickel, Paul Beardsley, Bob Sumner, and Markus Gross. 2010. High-Quality Single-Shot Capture of Facial Geometry. In ACM SIGGRAPH 2010 Papers (Los Angeles, California) (SIGGRAPH '10). Association for Computing Machinery, New York, NY, USA, Article 40, 9 pages. https://doi.org/10.1145/1833349.1778777
- Thabo Beeler, Markus Gross, Paulo Gotardo, Jérémy Riviere, and Derek Bradley. 2022. Medusa Facial Capture System. https://studios.disneyresearch.com/medusa
- Volker Blanz and Thomas Vetter. 1999. A morphable model for the synthesis of 3D faces. In Proceedings of the 26th annual conference on Computer graphics and interactive techniques. 187–194.
- Volker Blanz and Thomas Vetter. 2003. Face recognition based on fitting a 3D morphable model. IEEE Transactions on pattern analysis and machine intelligence 25, 9 (2003), 1063–1074.
- Blender. 2023. Rigify. https://docs.blender.org/manual/en/latest/addons/rigging/rigify/ introduction.html

- Timo Bolkart, Tianye Li, and Michael J Black. 2023. Instant Multi-View Head Capture through Learnable Registration. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 768–779.
- James Booth, Anastasios Roussos, Stefanos Zafeiriou, Allan Ponniah, and David Dunaway. 2016. A 3d morphable model learnt from 10,000 faces. In Proceedings of the IEEE conference on computer vision and pattern recognition. 5543–5552.
- Derek Bradley, Wolfgang Heidrich, Tiberiu Popa, and Alla Sheffer. 2010. High resolution passive facial performance capture. In ACM SIGGRAPH 2010 papers. 1–10.
- Adrian Bulat and Georgios Tzimiropoulos. 2017. 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. 1021–1030.
- Chen Cao, Tomas Simon, Jin Kyu Kim, Gabe Schwartz, Michael Zollhoefer, Shun-Suke Saito, Stephen Lombardi, Shih-En Wei, Danielle Belko, Shoou-I Yu, et al. 2022. Authentic volumetric avatars from a phone scan. ACM Transactions on Graphics (TOG) 41, 4 (2022), 1–19.
- Prashanth Chandran, Gaspard Zoss, Markus Gross, Paulo Gotardo, and Derek Bradley. 2022. Facial Animation with Disentangled Identity and Motion using Transformers. In ACM/Eurographics Symposium on Computer Animation 2022.
- Bindita Chaudhuri, Noranart Vesdapunt, Linda Shapiro, and Baoyuan Wang. 2020. Personalized Face Modeling for Improved Face Reconstruction and Motion Retargeting. In Computer Vision – ECCV 2020, Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm (Eds.). Springer International Publishing, Cham, 142–160.
- Grigorios G Chrysos, Epameinondas Antonakos, Patrick Snape, Akshay Asthana, and Stefanos Zafeiriou. 2018. A comprehensive performance evaluation of deformable face tracking "in-the-wild". *International Journal of Computer Vision* 126, 2 (2018), 198–232.
- Matthew Cong, Michael Bao, Jane L E, Kiran S Bhat, and Ronald Fedkiw. 2015. Fully automatic generation of anatomical face simulation models. In Proceedings of the 14th ACM SIGGRAPH/Eurographics Symposium on Computer Animation. 175–183.
- Daz3D. 2023. FaceGen. https://www.daz3d.com/easy-character-creation-with-facegen Paul Debevec. 2012. The light stages and their applications to photoreal digital actors. SIGGRAPH Asia Technical Briefs 2, 4 (2012), 1–6.
- Paul Debevec, Tim Hawkins, Chris Tchou, Haarm-Pieter Duiker, Westley Sarokin, and Mark Sagar. 2000. Acquiring the reflectance field of a human face. In Proceedings of the 27th annual conference on Computer graphics and interactive techniques. 145–156.
- Abdallah Dib, Junghyun Ahn, Cedric Thebault, Philippe-Henri Gosselin, and Louis Chevallier. 2023. S2f2: self-supervised high fidelity face reconstruction from monocular image. In 2023 IEEE 17th International Conference on Automatic Face and Gesture Recognition (FG). IEEE, 1–8.
- Akio Doi and Akio Koide. 1991. An efficient method of triangulating equi-valued surfaces by using tetrahedral cells. *IEICE TRANSACTIONS on Information and* Systems 74, 1 (1991), 214–224.
- Pengfei Dou and Ioannis A Kakadiaris. 2018. Multi-view 3D face reconstruction with deep recurrent neural networks. *Image and Vision Computing* 80 (2018), 80–91.
- P. Dou, S. K. Shah, and I. A. Kakadiaris. 2017. End-to-End 3D Face Reconstruction with Deep Neural Networks. In 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE Computer Society, Los Alamitos, CA, USA, 1503–1512.
- Epic Games. 2023. Unreal Engine. https://www.unrealengine.com
- Haiwen Feng, Timo Bolkart, Joachim Tesch, Michael J. Black, and Victoria Abrevaya. 2022. Towards Racially Unbiased Skin Tone Estimation via Scene Disambiguation. In European Conference on Computer Vision.
- David A Field. 1988. Laplacian smoothing and Delaunay triangulations. Communications in applied numerical methods 4, 6 (1988), 709–712.
- Andrean Franc. 2023. RIG LOGIC: RUNTIME EVALUATION OF METAHUMAN FACE RIGS. https://cdn2.unrealengine.com/rig-logic-whitepaper-v2-5c9f23f7e210.pdf. Accessed: 2023-05-30.
- Guy Gafni, Justus Thies, Michael Zollhöfer, and Matthias Nießner. 2021. Dynamic Neural Radiance Fields for Monocular 4D Facial Avatar Reconstruction. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). 8649–8658.
- Xuan Gao, Chenglai Zhong, Jun Xiang, Yang Hong, Yudong Guo, and Juyong Zhang. 2022. Reconstructing personalized semantic facial nerf models from monocular video. ACM Transactions on Graphics (TOG) 41, 6 (2022), 1–12.
- Ravi Garg, Anastasios Roussos, and Lourdes Agapito. 2013. Dense variational reconstruction of non-rigid surfaces from monocular video. In Proceedings of the IEEE Conference on computer vision and pattern recognition. 1272–1279.
- Pablo Garrido, Michael Zollhöfer, Dan Casas, Levi Valgaerts, Kiran Varanasi, Patrick Pérez, and Christian Theobalt. 2016. Reconstruction of personalized 3D face rigs from monocular video. ACM Transactions on Graphics (TOG) 35, 3 (2016), 1–15.
- Abhijeet Ghosh, Graham Fyffe, Borom Tunwattanapong, Jay Busch, Xueming Yu, and Paul Debevec. 2011. Multiview face capture using polarized spherical gradient illumination. In Proceedings of the 2011 SIGGRAPH Asia Conference. 1–10.
- John C Gower. 1975. Generalized procrustes analysis. Psychometrika 40 (1975), 33–51. Yudong Guo, Keyu Chen, Sen Liang, Yongjin Liu, Hujun Bao, and Juyong Zhang. 2021. AD-NeRF: Audio Driven Neural Radiance Fields for Talking Head Synthesis. In IEEE/CVF International Conference on Computer Vision (ICCV).

Yuxuan Han, Zhibo Wang, and Feng Xu. 2023. Learning a 3D Morphable Face Reflectance Model From Low-Cost Data. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 8598–8608.

HDReye Technologies Inc. 2023. HDReye. https://hdreye.app/

- Darren Hendler, Lucio Moser, Rishabh Battulwar, David Corral, Phil Cramer, Ron Miller, Rickey Cloudsdale, and Doug Roble. 2018. Avengers: Capturing Thanos's Complex Face. In ACM SIGGRAPH 2018 Talks (Vancouver, British Columbia, Canada) (SIGGRAPH '18). Association for Computing Machinery, New York, NY, USA, Article 58, 2 pages. https://doi.org/10.1145/3214745.3214766
- Dirk Holz, Alexandru E Ichim, Federico Tombari, Radu B Rusu, and Sven Behnke. 2015. Registration with the point cloud library: A modular framework for aligning in 3-D. *IEEE Robotics & Automation Magazine* 22, 4 (2015), 110–124.
- Hyprsense 2021. Hyprface Landmark Tracker. Retrieved Dec 23, 2023 from https: //medium.com/@hyprsense
- Alexandru Eugen Ichim, Sofien Bouaziz, and Mark Pauly. 2015. Dynamic 3D avatar creation from hand-held video input. ACM Transactions on Graphics (ToG) 34, 4 (2015), 1–14.
- Insta360. 2023. insta360 X3. https://store.insta360.com/product/x3
- Aaron S Jackson, Adrian Bulat, Vasileios Argyriou, and Georgios Tzimiropoulos. 2017. Large Pose 3D Face Reconstruction from a Single Image via Direct Volumetric CNN Regression. International Conference on Computer Vision (2017).
- Hanbyul Joo, Tomas Simon, Xulong Li, Hao Liu, Lei Tan, Lin Gui, Sean Banerjee, Timothy Scott Godisart, Bart Nabbe, Iain Matthews, Takeo Kanade, Shohei Nobuhara, and Yaser Sheikh. 2017. Panoptic Studio: A Massively Multiview System for Social Interaction Capture. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2017).
- Michael Kazhdan and Hugues Hoppe. 2013. Screened poisson surface reconstruction. ACM Transactions on Graphics (ToG) 32. 3 (2013). 1–13.
- Ira Kemelmacher-Shlizerman. 2013. Internet based morphable model. In Proceedings of the IEEE international conference on computer vision. 3256–3263.
- Ira Kemelmacher-Shlizerman and Steven M Seitz. 2011. Face reconstruction in the wild. In 2011 international conference on computer vision. IEEE, 1746–1753.
- Jongyoo Kim, Jiaolong Yang, and Xin Tong. 2021. Learning high-fidelity face texture completion without complete face texture. In Proceedings of the IEEE/CVF International Conference on Computer Vision. 13990–13999.
- Kurt Konolige. 1998. Small vision systems: Hardware and implementation. In *Robotics Research: The Eighth International Symposium.* Springer, 203–212.
- Alexandros Lattas, Yiming Lin, Jayanth Kannan, Ekin Ozturk, Luca Filipi, Giuseppe Claudio Guarnera, Gaurav Chawla, and Abhijeet Ghosh. 2022. Practical and Scalable Desktop-Based High-Quality Facial Capture. In Computer Vision – ECCV 2022. Springer Nature Switzerland, Cham, 522–537.
- John P Lewis, Ken Anjyo, Taehyun Rhee, Mengjie Zhang, Frederic H Pighin, and Zhigang Deng. 2014. Practice and theory of blendshape facial models. *Eurographics* (State of the Art Reports) 1, 8 (2014), 2.
- Moran Li, Haibin Huang, Yi Zheng, Mengtian Li, Nong Sang, and Chongyang Ma. 2022. Implicit Neural Deformation for Sparse-View Face Reconstruction. Computer Graphics Forum 41, 7 (2022), 601–610. https://doi.org/10.1111/cgf.14704
- Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier Romero. 2017. Learning a model of facial shape and expression from 4D scans. ACM Trans. Graph. 36, 6 (2017), 194–1.
- Shu Liang, Linda G Shapiro, and Ira Kemelmacher-Shlizerman. 2016. Head reconstruction from internet photos. In Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14. Springer, 360–374.
- Connor Lin, Koki Nagano, Jan Kautz, Eric Chan, Umar Iqbal, Leonidas Guibas, Gordon Wetzstein, and Sameh Khamis. 2023. Single-shot implicit morphable faces with consistent texture parameterization. In ACM SIGGRAPH 2023 Conference Proceedings. 1–12.
- Jiangke Lin, Yi Yuan, and Zhengxia Zou. 2021. Meingame: Create a game character face from a single portrait. In Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 35. 311–319.
- Winnie Lin, Yilin Zhu, Demi Guo, and Ron Fedkiw. 2022. Leveraging Deepfakes to Close the Domain Gap between Real and Synthetic Images in Facial Capture Pipelines. arXiv:2204.10746 [cs.CV]
- Shichen Liu, Yunxuan Cai, Haiwei Chen, Yichao Zhou, and Yajie Zhao. 2022. Rapid Face Asset Acquisition with Recurrent Feature Alignment. ACM Trans. Graph. 41, 6, Article 214 (nov 2022), 17 pages.
- M2M 2023. Mesh to Metahuman for Unreal Engine. Retrieved Dec 23, 2023 from https://dev.epicgames.com/documentation/en-us/metahuman/mesh-tometahuman-for-unreal-engine
- MHC 2023. Metahuman. Retrieved Dec 23, 2023 from https://dev.epicgames.com/ documentation/en-US/metahuman/metahuman-documentation
- Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Commun. ACM* 65, 1 (2021), 99–106.

- Marcel Piotraschke and Volker Blanz. 2016. Automated 3d face reconstruction from multiple images using quality measures. In *Proceedings of the IEEE conference on computer vision and pattern recognition.* 3418–3427.
- Dafei Qin, Jun Saito, Noam Aigerman, Thibault Groueix, and Taku Komura. 2023. Neural Face Rigging for Animating and Retargeting Facial Meshes in the Wild. In ACM SIGGRAPH 2023 Conference Proceedings (Los Angeles, CA, USA) (SIGGRAPH '23). Association for Computing Machinery, New York, NY, USA, Article 68, 11 pages. https://doi.org/10.1145/3588432.3591556
- Gilles Rainer, Lewis Bridgeman, and Abhijeet Ghosh. 2023. Neural shading fields for efficient facial inverse rendering. In *Computer Graphics Forum*, Vol. 42. Wiley Online Library, e14943.
- Ravi Ramamoorthi and Pat Hanrahan. 2001. An Efficient Representation for Irradiance Environment Maps. In Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH '01). Association for Computing Machinery, New York, NY, USA, 497–500.
- Nick Rasmussen, Duc Quang Nguyen, Willi Geiger, and Ronald Fedkiw. 2003. Smoke simulation for large scale phenomena. In ACM SIGGRAPH 2003 Papers. 703–707.
- Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. 2020. Accelerating 3D Deep Learning with Py-Torch3D. arXiv:2007.08501 (2020).
- Reallusion. 2023. Headshot2. https://www.reallusion.com/character-creator/headshot/ Xingyu Ren, Jiankang Deng, Chao Ma, Yichao Yan, and Xiaokang Yang. 2023. Improving Fairness in Facial Albedo Estimation via Visual-Textual Cues. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 4511–4520.
- E. Richardson, M. Sela, and R. Kimmel. 2016. 3D Face Reconstruction by Learning from Synthetic Data. In 2016 Fourth International Conference on 3D Vision (3DV). IEEE Computer Society, Los Alamitos, CA, USA, 460–469.
- RICOH360. 2023. RICOH THETA. https://www.ricoh360.com/theta/
- Jérémy Riviere, Paulo Gotardo, Derek Bradley, Abhijeet Ghosh, and Thabo Beeler. 2020. Single-Shot High-Quality Facial Geometry and Skin Appearance Capture. ACM Trans. Graph. 39, 4, Article 81 (jul 2020), 12 pages. https://doi.org/10.1145/3386569. 3392464
- Romdhani and Vetter. 2003. Efficient, robust and accurate fitting of a 3D morphable model. In Proceedings Ninth IEEE International Conference on Computer Vision. 59–66 vol.1. https://doi.org/10.1109/ICCV.2003.1238314
- Joseph Roth, Yiying Tong, and Xiaoming Liu. 2016. Adaptive 3D face reconstruction from unconstrained photo collections. In Proceedings of the IEEE conference on computer vision and pattern recognition. 4197–4206.
- Michael Sailer, Jan Ulrich Hense, Sarah Katharina Mayr, and Heinz Mandl. 2017. How gamification motivates: An experimental study of the effects of specific game design elements on psychological need satisfaction. *Computers in human behavior* 69 (2017), 371–380.
- Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. 2019. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In Proceedings of the IEEE/CVF international conference on computer vision. 2304–2314.
- Soubhik Sanyal, Timo Bolkart, Haiwen Feng, and Michael J Black. 2019. 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. 7763–7772.
- Soumyadip Sengupta, Brian Curless, Ira Kemelmacher-Shlizerman, and Steven M Seitz. 2021. A light stage on every desk. In Proceedings of the IEEE/CVF International Conference on Computer Vision. 2420–2429.
- Soumyadip Sengupta, Angjoo Kanazawa, Carlos D. Castillo, and David W. Jacobs. 2018. SfSNet: Learning Shape, Refectance and Illuminance of Faces in the Wild. In *Computer Vision and Pattern Regognition (CVPR).*
- James A Sethian. 1999. Fast marching methods. SIAM review 41, 2 (1999), 199-235.
- Artem Sevastopolsky, Savva Ignatiev, Gonzalo Ferrer, Evgeny Burnaev, and Victor Lempitsky. 2020. Relightable 3d head portraits from a smartphone video. arXiv preprint arXiv:2012.09963 (2020).
- Fuhao Shi, Hsiang-Tao Wu, Xin Tong, and Jinxiang Chai. 2014. Automatic acquisition of high-fidelity facial performances using monocular videos. *ACM Transactions on Graphics (TOG)* 33, 6 (2014), 1–13.
- Tianyang Shi, Zhengxia Zuo, Yi Yuan, and Changjie Fan. 2020. Fast and robust face-toparameter translation for game character auto-creation. In Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 34. 1733–1740.
- Ron Slossberg, Ibrahim Jubran, and Ron Kimmel. 2022. Unsupervised High-Fidelity Facial Texture Generation and Reconstruction. In *Computer Vision – ECCV 2022*, Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner (Eds.). Springer Nature Switzerland, Cham, 212–229.
- William AP Smith, Alassane Seck, Hannah Dee, Bernard Tiddeman, Joshua B Tenenbaum, and Bernhard Egger. 2020. A morphable face albedo model. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 5011–5020.
- Ayush Tewari, Florian Bernard, Pablo Garrido, Gaurav Bharaj, Mohamed Elgharib, Hans-Peter Seidel, Patrick Pérez, Michael Zöllhofer, and Christian Theobalt. 2019. Fml: Face model learning from videos. In Proceedings of the IEEE Conference on

16 · Zhu, Y. et al

Computer Vision and Pattern Recognition. 10812–10822.

- A. Tewari, J. Thies, B. Mildenhall, P. Srinivasan, E. Tretschk, W. Yifan, C. Lassner, V. Sitzmann, R. Martin-Brualla, S. Lombardi, T. Simon, C. Theobalt, M. Nießner, J. T. Barron, G. Wetzstein, M. Zollhöfer, and V. Golyanik. 2022. Advances in Neural Rendering. *Computer Graphics Forum* 41, 2 (2022), 703–735. _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1111/cgf.14507.
- Ayush Tewari, Michael Zollhoefer, Florian Bernard, Pablo Garrido, Hyeongwoo Kim, Patrick Perez, and Christian Theobalt. 2018. High-fidelity monocular face reconstruction based on an unsupervised model-based face autoencoder. *IEEE transactions* on pattern analysis and machine intelligence 42, 2 (2018), 357–370.
- Ayush Tewari, Michael Zollhöfer, Hyeongwoo Kim, Pablo Garrido, Florian Bernard, Patrick Pérez, and Christian Theobalt. 2017. MoFA: Model-Based Deep Convolutional Face Autoencoder for Unsupervised Monocular Reconstruction. In 2017 IEEE International Conference on Computer Vision (ICCV). 3735–3744.
- The Blender Foundation. 2023. Blender. https://www.blender.org/
- Luan Tran, Feng Liu, and Xiaoming Liu. 2019. Towards High-Fidelity Nonlinear 3D Face Morphable Model. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).
- L. Tran and X. Liu. 2018. Nonlinear 3D Face Morphable Model. In 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). IEEE Computer Society, Los Alamitos, CA, USA, 7346–7355.
- Edith Tretschk, Navami Kairanda, Mallikarjun BR, Rishabh Dabral, Adam Kortylewski, Bernhard Egger, Marc Habermann, Pascal Fua, Christian Theobalt, and Vladislav Golyanik. 2023. State of the Art in Dense Monocular Non-Rigid 3D Reconstruction. In *Computer Graphics Forum*, Vol. 42. Wiley Online Library, 485–520.
- Lizhen Wang, Zhiyuan Chen, Tao Yu, Chenguang Ma, Liang Li, and Yebin Liu. 2022. Faceverse: a fine-grained and detail-controllable 3d face morphable model from a hybrid dataset. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 20333-20342.
- Yifan Wang, Aleksander Holynski, Xiuming Zhang, and Xuaner Zhang. 2023. Sunstage: Portrait reconstruction and relighting using the sun as a light stage. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 20792–20802.
- Ziyan Wang, Timur Bagautdinov, Stephen Lombardi, Tomas Simon, Jason Saragih, Jessica Hodgins, and Michael Zollhofer. 2021. Learning Compositional Radiance Fields of Dynamic Human Heads. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). 5704–5713.
- Antony Ward. 2004. Game character development with maya. New Riders.
- F. Wu, L. Bao, Y. Chen, Y. Ling, Y. Song, S. Li, K. Ngan, and W. Liu. 2019. MVF-Net: Multi-View 3D Face Morphable Model Regression. In 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). IEEE Computer Society, Los Alamitos, CA, USA, 959–968.
- Jane Wu, Michael Bao, Xinwei Yao, and Ronald Fedkiw. 2023a. Deep Energies for Estimating Three-Dimensional Facial Pose and Expression. Communications on Applied Mathematics and Computation (2023), 1–25.
- Yiqian Wu, Jing Zhang, Hongbo Fu, and Xiaogang Jin. 2023b. LPFF: A Portrait Dataset for Face Generators Across Large Poses. arXiv preprint arXiv:2303.14407 (2023).
- Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. 2020. Multiview Neural Surface Reconstruction by Disentangling Geometry and Appearance. In Advances in Neural Information Processing Systems, H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 2492–2502.
- Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. 2018. BiSeNet: Bilateral Segmentation Network for Real-time Semantic Segmentation. In Proceedings of the European Conference on Computer Vision (ECCV).
- Xiaoxing Zeng, Xiaojiang Peng, and Yu Qiao. 2019. DF2Net: A Dense-Fine-Finer Network for Detailed 3D Face Reconstruction. In Proceedings of the IEEE International Conference on Computer Vision. 2315–2324.
- Jingbo Zhang, Xiaoyu Li, Ziyu Wan, Can Wang, and Jing Liao. 2022b. Fdnerf: Few-shot dynamic neural radiance fields for face reconstruction and expression editing. In SIGGRAPH Asia 2022 Conference Papers. 1–9.
- Li Zhang, Noah Snavely, Brian Curless, and Steven M. Seitz. 2004. Spacetime Faces: High-Resolution Capture for Modeling and Animation. In ACM Annual Conference on Computer Graphics (Los Angeles, CA). 548–558.
- Longwen Zhang, Chuxiao Zeng, Qixuan Zhang, Hongyang Lin, Ruixiang Cao, Wei Yang, Lan Xu, and Jingyi Yu. 2022c. Video-driven neural physically-based facial asset for production. ACM Transactions on Graphics (TOG) 41, 6 (2022), 1–16.
- Zhenyu Zhang, Yanhao Ge, Ying Tai, Weijian Cao, Renwang Chen, Kunlin Liu, Hao Tang, Xiaoming Huang, Chengjie Wang, Zhifeng Xie, et al. 2022a. Physicallyguided Disentangled Implicit Rendering for 3D Face Modeling. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 20353–20363.
- Mingwu Zheng, Haiyu Zhang, Hongyu Yang, and Di Huang. 2023b. NeuFace: Realistic 3D Neural Face Rendering from Multi-view Images. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 16868–16877.
- Yufeng Zheng, Wang Yifan, Gordon Wetzstein, Michael J Black, and Otmar Hilliges. 2023a. Pointavatar: Deformable point-based head avatars from videos. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 21057–21067.

- X. Zhu, Z. Lei, X. Liu, H. Shi, and S. Z. Li. 2016. Face Alignment Across Large Poses: A 3D Solution. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE Computer Society, Los Alamitos, CA, USA, 146–155.
- Michael Zollhöfer, Justus Thies, Pablo Garrido, Derek Bradley, Thabo Beeler, Patrick Pérez, Marc Stamminger, Matthias Nießner, and Christian Theobalt. 2018. State of the art on monocular 3D face reconstruction, tracking, and applications. In *Computer Graphics Forum*, Vol. 37. Wiley Online Library, 523–550.